
Scenario Reduction Revisited: Fundamental Limits and Guarantees

Napat Rujeerapaiboon, Kilian Schindler,
Daniel Kuhn, Wolfram Wiesemann

Abstract The goal of scenario reduction is to approximate a given discrete distribution with another discrete distribution that has fewer atoms. We distinguish continuous scenario reduction, where the new atoms may be chosen freely, and discrete scenario reduction, where the new atoms must be chosen from among the existing ones. Using the Wasserstein distance as measure of proximity between distributions, we identify those n -point distributions on the unit ball that are least susceptible to scenario reduction, *i.e.*, that have maximum Wasserstein distance to their closest m -point distributions for some prescribed $m < n$. We also provide sharp bounds on the added benefit of continuous over discrete scenario reduction. Finally, to our best knowledge, we propose the first polynomial-time constant-factor approximations for both discrete and continuous scenario reduction as well as the first exact exponential-time algorithms for continuous scenario reduction.

Keywords scenario reduction, Wasserstein distance, constant-factor approximation algorithm, k -median clustering, k -means clustering

1 Introduction

The vast majority of numerical solution schemes in stochastic programming rely on a discrete approximation of the true (typically continuous) probability distribution governing the uncertain problem parameters. This discrete approximation is often generated by sampling from the true distribution. Alternatively, it could be constructed directly from real historical observations of the uncertain parameters.

Napat Rujeerapaiboon, Kilian Schindler, Daniel Kuhn
Risk Analytics and Optimization Chair
École Polytechnique Fédérale de Lausanne, Switzerland
Tel.: +41 (0)21 693 00 36 Fax: +41 (0)21 693 24 89
E-mail: napat.rujeerapaiboon@epfl.ch, kilian.schindler@epfl.ch, daniel.kuhn@epfl.ch

Wolfram Wiesemann
Imperial College Business School
Imperial College London, United Kingdom
Tel.: +44 (0)20 7594 9150
E-mail: ww@imperial.ac.uk

To obtain a faithful approximation for the true distribution, however, the discrete distribution must have a large number n of support points or *scenarios*, which may render the underlying stochastic program computationally excruciating.

An effective means to ease the computational burden is to rely on *scenario reduction* pioneered by Dupačová et al (2003), which aims to approximate the initial n -point distribution with a simpler m -point distribution ($m < n$) that is as close as possible to the initial distribution with respect to a probability metric; see also Heitsch and Römisch (2003). The modern stability theory of stochastic programming surveyed by Dupačová (1990) and Römisch (2003) indicates that the Wasserstein distance may serve as a natural candidate for this probability metric.

Our interest in Wasserstein distance-based scenario reduction is also fuelled by recent progress in data-driven distributionally robust optimization, where it has been shown that the worst-case expectation of an uncertain cost over all distributions in a Wasserstein ball can often be computed efficiently via convex optimization (Mohajerin Esfahani and Kuhn 2015, Zhao and Guan 2015, Gao and Kleywegt 2016). A Wasserstein ball is defined as the family of all distributions that are within a certain Wasserstein distance from a discrete reference distribution. As distributionally robust optimization problems over Wasserstein balls are harder to solve than their stochastic counterparts, we expect significant computational savings from replacing the initial n -point reference distribution with a new m -point reference distribution. The benefits of scenario reduction may be particularly striking for two-stage distributionally robust linear programs, which admit tight approximations as semidefinite programs (Hanasusanto and Kuhn 2016).

Suppose now that the initial distribution is given by $\mathbb{P} = \sum_{i \in I} p_i \delta_{\xi_i}$, where $\xi_i \in \mathbb{R}^d$ and $p_i \in [0, 1]$ represent the location and probability of the i -th scenario of \mathbb{P} for $i \in I = \{1, \dots, n\}$. Similarly, assume that the reduced target distribution is representable as $\mathbb{Q} = \sum_{j \in J} q_j \delta_{\zeta_j}$, where $\zeta_j \in \mathbb{R}^d$ and $q_j \in [0, 1]$ stand for the location and probability of the j -th scenario of \mathbb{Q} for $j \in J = \{1, \dots, m\}$. Then, the type- l Wasserstein distance between \mathbb{P} and \mathbb{Q} is defined through

$$d_l(\mathbb{P}, \mathbb{Q}) = \left[\min_{\Pi \in \mathbb{R}_+^{n \times m}} \left\{ \sum_{i \in I} \sum_{j \in J} \pi_{ij} \|\xi_i - \zeta_j\|^l : \begin{array}{l} \sum_{j \in J} \pi_{ij} = p_i \quad \forall i \in I \\ \sum_{i \in I} \pi_{ij} = q_j \quad \forall j \in J \end{array} \right\} \right]^{1/l},$$

where $l \geq 1$ and $\|\cdot\|$ denotes some norm on \mathbb{R}^d , see, *e.g.*, Heitsch and Römisch (2007) or Pflug and Pichler (2011). The linear program in the definition of the Wasserstein distance can be viewed as a minimum-cost transportation problem, where π_{ij} represents the amount of probability mass shipped from ξ_i to ζ_j at unit transportation cost $\|\xi_i - \zeta_j\|^l$. Thus, $d_l^l(\mathbb{P}, \mathbb{Q})$ quantifies the minimum cost of moving the initial distribution \mathbb{P} to the target distribution \mathbb{Q} .

For any $\Xi \subseteq \mathbb{R}^d$, we denote by $\mathcal{P}_{\Xi}(\Xi, n)$ the set of all uniform discrete distributions on Ξ with exactly n distinct scenarios and by $\mathcal{P}(\Xi, m)$ the set of all (not necessarily uniform) discrete distributions on Ξ with at most m scenarios. We henceforth assume that $\mathbb{P} \in \mathcal{P}_{\Xi}(\mathbb{R}^d, n)$. This assumption is crucial for the simplicity of the results in Sections 2 and 3, and it is almost surely satisfied whenever \mathbb{P} is obtained via sampling from a continuous probability distribution. Hence, we can think of \mathbb{P} as an *empirical distribution*. To remind us of this interpretation, we will henceforth denote the initial distribution by $\hat{\mathbb{P}}_n$. Note that the pairwise difference of the scenarios can always be enforced by slightly perturbing their locations, while

the uniformity of their probabilities can be enforced by decomposing the scenarios into clusters of close but mutually distinct sub-scenarios with (smaller) uniform probabilities.

We are now ready to introduce the *continuous scenario reduction problem*

$$C_l(\hat{\mathbb{P}}_n, m) = \min_{\mathbb{Q}} \left\{ d_l(\hat{\mathbb{P}}_n, \mathbb{Q}) : \mathbb{Q} \in \mathcal{P}(\mathbb{R}^d, m) \right\},$$

where the new scenarios ζ_j , $j \in J$, of the target distribution \mathbb{Q} may be chosen freely from within \mathbb{R}^d , as well as the *discrete scenario reduction problem*

$$D_l(\hat{\mathbb{P}}_n, m) = \min_{\mathbb{Q}} \left\{ d_l(\hat{\mathbb{P}}_n, \mathbb{Q}) : \mathbb{Q} \in \mathcal{P}(\text{supp}(\hat{\mathbb{P}}_n), m) \right\},$$

where the new scenarios must be chosen from within the support of the empirical distribution, which is given by the finite set $\text{supp}(\hat{\mathbb{P}}_n) = \{\xi_i : i \in I\}$. Even though the continuous scenario reduction problem offers more flexibility and is therefore guaranteed to find (weakly) better approximations to the initial empirical distribution, to our best knowledge, the existing stochastic programming literature has exclusively focused on the discrete scenario reduction problem.

Note that if the support points ζ_j , $j \in J$, are fixed, then both scenario reduction problems simplify to a linear program over the probabilities q_j , $j \in J$, which admits an explicit solution (Dupačová et al 2003, Theorem 2). Otherwise, however, both problems are intractable. Indeed, if $l = 1$, then the discrete scenario reduction problem represents a metric k -median problem with $k = m$, which was shown to be \mathcal{NP} -hard by Kariv and Hakimi (1979). If $l = 2$ and distances in \mathbb{R}^d are measured by the 2-norm, on the other hand, then the continuous scenario reduction problem constitutes a k -means clustering problem with $k = m$, which is \mathcal{NP} -hard even if $d = 2$ or $m = 2$; see Mahajan et al (2009) and Aloise et al (2009).

Heitsch and Römisch (2003) have shown that the discrete scenario reduction problem admits a reformulation as a mixed-integer linear program (MILP), which can be solved to global optimality for $n \lesssim 10^3$ using off-the-shelf solvers. For larger instances, however, one must resort to approximation algorithms. Most large-scale discrete scenario reduction problems are nowadays solved with a greedy heuristic that was originally devised by Dupačová et al (2003) and further refined by Heitsch and Römisch (2003). For example, this heuristic is routinely used for scenario (tree) reduction in the context of power systems operations, see, *e.g.*, Römisch and Vigerske (2010) or Morales et al (2009) and the references therein. Despite its practical success, we will show in Section 4 that this heuristic fails to provide a constant-factor approximation for the discrete scenario reduction problem.

This paper extends the theory of scenario reduction along several dimensions.

- (i) We establish fundamental performance guarantees for continuous scenario reduction when $l \in \{1, 2\}$, *i.e.*, we show that the Wasserstein distance of the initial n -point distribution to its nearest m -point distribution is bounded by $\sqrt{\frac{n-m}{n-1}}$ across all initial distributions on the unit ball in \mathbb{R}^d . We show that for $l = 2$ this worst-case performance is attained by some initial distribution, which we construct explicitly. We also provide evidence indicating that this worst-case performance reflects the norm rather than the exception in high dimensions d . Finally, we provide a lower bound on the worst-case performance for $l = 1$.

- (ii) We analyze the loss of optimality incurred by solving the discrete scenario reduction problem instead of its continuous counterpart. Specifically, we demonstrate that the ratio $D_l(\hat{\mathbb{P}}_n, m)/C_l(\hat{\mathbb{P}}_n, m)$ is bounded by $\sqrt{2}$ for $l = 2$ and by 2 for $l = 1$. We also show that these bounds are essentially tight.
- (iii) We showcase the intimate relation between scenario reduction and k -means clustering. By leveraging existing constant-factor approximation algorithms for k -median clustering problems due to Arya et al (2004) and the new performance bounds from (ii), we develop the first polynomial-time constant-factor approximation algorithms for both continuous and discrete scenario reduction. We also show that these algorithms can be warmstarted using the greedy heuristic by Dupačová et al (2003) to improve practical performance.
- (iv) We present exact mixed-integer programming reformulations for the continuous scenario reduction problem.

Continuous scenario reduction is intimately related to the optimal quantization of probability distributions, where one seeks an m -point distribution approximating a non-discrete initial distribution. Research efforts in this domain have mainly focused on the asymptotic behavior of the quantization problem as m tends to infinity, see Graf and Luschgy (2000). The ramifications of this stream of literature for stochastic programming are discussed by Pflug and Pichler (2011). Techniques familiar from scenario reduction lend themselves also for scenario generation, where one aims to construct a scenario tree with a prescribed branching structure that approximates a given stochastic process with respect to a probability metric, see, e.g., Pflug (2001) and Hochreiter and Pflug (2007).

The rest of this paper unfolds as follows. Section 2 seeks to identify n -point distributions on the unit ball that are least susceptible to scenario reduction, *i.e.*, that have maximum Wasserstein distance to their closest m -point distributions, and Section 3 discusses sharp bounds on the added benefit of continuous over discrete scenario reduction. Section 4 presents exact exponential-time algorithms as well as polynomial-time constant-factor approximations for scenario reduction. Section 5 reports on numerical results for a color quantization experiment. Unless otherwise specified, below we will always work with the 2-norm on \mathbb{R}^d .

Notation: We let \mathbb{I} be the identity matrix, \mathbf{e} the vector of all ones and \mathbf{e}_i the i -th standard basis vector of appropriate dimensions. The ij -th element of a matrix \mathbf{A} is denoted by a_{ij} . For \mathbf{A} and \mathbf{B} in the space \mathbb{S}^n of symmetric $n \times n$ matrices, the relation $\mathbf{A} \succeq \mathbf{B}$ means that $\mathbf{A} - \mathbf{B}$ is positive semidefinite. Generic norms are denoted by $\|\cdot\|$, while $\|\cdot\|_p$ stands for the p -norm, $p \geq 1$. For $\Xi \subseteq \mathbb{R}^d$, we define $\mathcal{P}(\Xi, m)$ as the set of all probability distributions supported on at most m points in Ξ and $\mathcal{P}_{\text{E}}(\Xi, n)$ as the set of all *uniform* distributions supported on exactly n *distinct* points in Ξ . The support of a probability distribution \mathbb{P} is denoted by $\text{supp}(\mathbb{P})$, and the Dirac distribution concentrating unit mass at $\boldsymbol{\xi}$ is denoted by $\delta_{\boldsymbol{\xi}}$.

2 Fundamental Limits of Scenario Reduction

In this section we characterize the Wasserstein distance $C_l(\hat{\mathbb{P}}_n, m)$ between an n -point empirical distribution $\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\boldsymbol{\xi}_i}$ and its continuously reduced optimal m -point distribution $\mathbb{Q} \in \mathcal{P}(\mathbb{R}^d, m)$. Since the positive homogeneity of the Wasserstein distance d_l implies that $C_l(\hat{\mathbb{P}}_n^{\lambda}, m) = \lambda \cdot C_l(\hat{\mathbb{P}}_n, m)$ for the scaled distribution

$\hat{\mathbb{P}}'_n = \frac{1}{n} \sum_{i=1}^n \delta_{\lambda \xi_i}$, $\lambda \in \mathbb{R}_+$, we restrict ourselves to empirical distributions $\hat{\mathbb{P}}_n$ whose scenarios satisfy $\|\xi_i\|_2 \leq 1$, $i = 1, \dots, n$. We thus want to quantify

$$\overline{C}_l(n, m) = \max_{\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)} \left\{ C_l(\hat{\mathbb{P}}_n, m) : \|\xi\|_2 \leq 1 \ \forall \xi \in \text{supp}(\hat{\mathbb{P}}_n) \right\}, \quad (1)$$

which amounts to the worst-case (*i.e.*, largest) Wasserstein distance between any n -point empirical distribution $\hat{\mathbb{P}}_n$ over the unit ball and its optimally selected continuous m -point scenario reduction. By construction, this worst-case distance satisfies $\overline{C}_l(n, m) \geq 0$, and the lower bound is attained whenever $n = m$. One also verifies that $\overline{C}_l(n, m) \leq \overline{C}_l(n, 1) \leq 1$ since the Wasserstein distance to the Dirac distribution δ_0 is bounded above by 1. Our goal is to derive possibly tight upper bounds on $\overline{C}_l(n, m)$ for the Wasserstein distances of type $l \in \{1, 2\}$.

In the following, we denote by $\mathfrak{P}(I, m)$ the family of all m -set partitions of the index set I , *i.e.*,

$$\mathfrak{P}(I, m) = \left\{ \{I_1, \dots, I_m\} : \emptyset \neq I_1, \dots, I_m \subseteq I, \cup_j I_j = I, I_i \cap I_j = \emptyset \ \forall i \neq j \right\},$$

and an element of this set (*i.e.* a specific m -set partition) as $\{I_j\} \in \mathfrak{P}(I, m)$. Our derivations will make extensive use of the following theorem.

Theorem 1 *For any type- l Wasserstein distance induced by any norm $\|\cdot\|$, the continuous scenario reduction problem can be reformulated as*

$$C_l(\hat{\mathbb{P}}_n, m) = \min_{\{I_j\} \in \mathfrak{P}(I, m)} \left[\frac{1}{n} \sum_{j \in J} \min_{\zeta_j \in \mathbb{R}^d} \sum_{i \in I_j} \|\xi_i - \zeta_j\|^l \right]^{1/l}. \quad (2)$$

Problem (2) can be interpreted as a Voronoi partitioning problem that asks for a Voronoi decomposition of \mathbb{R}^d into m cells whose Voronoi centroids ζ_1, \dots, ζ_m minimize the cumulative l -th powers of the distances to n prespecified points ξ_1, \dots, ξ_n .

Proof of Theorem 1 Theorem 2 of Dupačová et al (2003) implies that the smallest Wasserstein distance between the empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ and any distribution \mathbb{Q} supported on a finite set $\Xi \subset \mathbb{R}^d$ amounts to

$$\min_{\mathbb{Q} \in \mathcal{P}(\Xi, \infty)} d_l(\hat{\mathbb{P}}_n, \mathbb{Q}) = \left[\frac{1}{n} \sum_{i \in I} \min_{\zeta \in \Xi} \|\xi_i - \zeta\|^l \right]^{1/l},$$

where $\mathcal{P}(\Xi, \infty)$ denotes the set of all probability distributions supported on the finite set Ξ . The continuous scenario reduction problem $C_l(\hat{\mathbb{P}}_n, m)$ selects the set Ξ^* that minimizes this quantity over all sets in $\Xi \subset \mathbb{R}^d$ with $|\Xi| = m$ elements:

$$C_l(\hat{\mathbb{P}}_n, m) = \min_{\{\zeta_j\} \subseteq \mathbb{R}^d} \left[\frac{1}{n} \sum_{i \in I} \min_{j \in J} \|\xi_i - \zeta_j\|^l \right]^{1/l}. \quad (3)$$

One readily verifies that any optimal solution $\{\zeta_1^*, \dots, \zeta_m^*\}$ to problem (3) corresponds to an optimal solution $\{I_1^*, \dots, I_m^*\}$ to problem (2) with the same objective value if we identify the set I_j^* with all observations ξ_i that are closer to ζ_j^* than any other $\zeta_{j'}^*$ (ties may be broken arbitrarily). Likewise, any optimal solution $\{I_1^*, \dots, I_m^*\}$ to problem (2) with inner minimizers $\{\zeta_1^*, \dots, \zeta_m^*\}$ translates into an optimal solution $\{\zeta_1^*, \dots, \zeta_m^*\}$ to problem (3) with the same objective value. \square

Remark 1 (Minimizers of (2)) For $l = 2$, the inner minimum corresponding to the set I_j is attained by the *mean* $\zeta_j^* = \text{mean}(I_j) = \frac{1}{|I_j|} \sum_{i \in I_j} \xi_i$. Likewise, for $l = 1$, the inner minimum corresponding to the set I_j is attained by any *geometric median*

$$\zeta_j^* = \text{gmed}(I_j) \in \arg \min_{\zeta_j \in \mathbb{R}^d} \sum_{i \in I_j} \|\xi_i - \zeta_j\|,$$

which can be determined efficiently by solving a second-order cone program whenever a p -norm with rational $p \geq 1$ is considered (Alizadeh and Goldfarb (2003)).

The rest of this section derives tight upper bounds on $\bar{C}_l(n, m)$ for Wasserstein distances of type $l = 2$ (Section 2.1) as well as upper and lower bounds for Wasserstein distances of type $l = 1$ (Section 2.2). We summarize and discuss our findings in Section 2.3.

2.1 Fundamental Limits for the Type-2 Wasserstein Distance

We now derive a revised upper bound on $\bar{C}_l(n, m)$ for the type-2 Wasserstein distance. The result relies on auxiliary lemmas that are relegated to the appendix.

Theorem 2 *The worst-case type-2 Wasserstein distance satisfies $\bar{C}_2(n, m) \leq \sqrt{\frac{n-m}{n-1}}$.*

Note that whenever the reduced distribution satisfies $m > 1$, the bound of Theorem 2 is strictly tighter than the naïve bound of 1 from the previous section.

Proof of Theorem 2 From Theorem 1 and Remark 1 we observe that

$$\begin{aligned} \bar{C}_2(n, m) &= \max_{\{\xi_i\} \subseteq \mathbb{R}^d} \min_{\{I_j\} \in \mathfrak{P}(I, m)} \left[\frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2^2 \right]^{1/2} \\ \text{s.t.} \quad &\|\xi_i\|_2 \leq 1 \quad \forall i \in I. \end{aligned}$$

Introducing the epigraphical variable τ , this problem can be expressed as

$$\begin{aligned} \bar{C}_2^2(n, m) &= \max_{\tau \in \mathbb{R}, \{\xi_i\} \subseteq \mathbb{R}^d} \frac{1}{n} \tau \\ \text{s.t.} \quad &\tau \leq \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2^2 \quad \forall \{I_j\} \in \mathfrak{P}(I, m) \\ &\xi_i^\top \xi_i \leq 1 \quad \forall i \in I. \end{aligned} \quad (4)$$

For each $j \in J$ and $i \in I_j$, the squared norm in the first constraint of (4) can be expressed in terms of the inner products between pairs of empirical observations:

$$\begin{aligned} \|\xi_i - \text{mean}(I_j)\|_2^2 &= \frac{1}{|I_j|^2} \left\| |I_j| \xi_i - \sum_{k \in I_j} \xi_k \right\|_2^2 \\ &= \frac{1}{|I_j|^2} \left(|I_j|^2 \xi_i^\top \xi_i - 2|I_j| \sum_{k \in I_j} \xi_i^\top \xi_k + \sum_{k \in I_j} \xi_k^\top \xi_k + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} \xi_k^\top \xi_{k'} \right). \end{aligned}$$

Introducing the Gram matrix

$$\mathbf{S} = [\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n]^\top [\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n] \in \mathbb{S}^n, \quad \mathbf{S} \succeq \mathbf{0} \text{ and } \text{rank}(\mathbf{S}) \leq \min\{n, d\} \quad (5)$$

then allows us to simplify the first constraint in (4) to

$$\tau \leq \sum_{j \in J} \frac{1}{|I_j|^2} \sum_{i \in I_j} \left(|I_j|^2 s_{ii} - 2|I_j| \sum_{k \in I_j} s_{ik} + \sum_{k \in I_j} s_{kk} + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} s_{kk'} \right).$$

Note that the second constraint in problem (4) can now be expressed as $s_{ii} \leq 1$, and hence all constraints in (4) are linear in the Gram matrix \mathbf{S} .

Our discussion implies that we obtain an upper bound on $\bar{\mathcal{C}}_2(n, m)$ by reformulating problem (4) as a semidefinite program in terms of the Gram matrix \mathbf{S}

$$\begin{aligned} & \max_{\tau \in \mathbb{R}, \mathbf{S} \in \mathbb{S}^n} \quad \frac{1}{n} \tau \\ & \text{s.t.} \quad \tau \leq \sum_{j \in J} \frac{1}{|I_j|^2} \sum_{i \in I_j} \left(|I_j|^2 s_{ii} - 2|I_j| \sum_{k \in I_j} s_{ik} + \sum_{k \in I_j} s_{kk} + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} s_{kk'} \right) \\ & \quad \quad \quad \forall \{I_j\} \in \mathfrak{P}(I, m) \\ & \quad \quad \mathbf{S} \succeq \mathbf{0}, \quad s_{ii} \leq 1 \quad \forall i \in I, \end{aligned} \quad (6)$$

where we have relaxed the rank condition in the definition of the Gram matrix (5). Lemma 1 in the appendix shows that (6) has an optimal solution (τ^*, \mathbf{S}^*) that satisfies $\mathbf{S}^* = \alpha \mathbb{I} + \beta \mathbf{e}\mathbf{e}^\top$ for some $\alpha, \beta \in \mathbb{R}$. Moreover, Lemma 2 in the appendix shows that any matrix of the form $\mathbf{S} = \alpha \mathbb{I} + \beta \mathbf{1}\mathbf{1}^\top$ is positive semidefinite if and only if $\alpha \geq 0$ and $\alpha + n\beta \geq 0$. We thus conclude that (6) can be reformulated as

$$\begin{aligned} & \max_{\tau, \alpha, \beta \in \mathbb{R}} \quad \frac{1}{n} \tau \\ & \text{s.t.} \quad \tau \leq (n - m)\alpha, \quad \alpha + \beta \leq 1 \\ & \quad \quad \alpha \geq 0, \quad \alpha + n\beta \geq 0, \end{aligned} \quad (7)$$

where the first constraint follows from the fact that for any set I_j in (6), we have

$$\frac{1}{|I_j|^2} \sum_{i \in I_j} \left(|I_j|^2 (\alpha + \beta) - 2|I_j| (\alpha + |I_j|\beta) + \sum_{k \in I_j} (\alpha + \beta) + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} \beta \right) = (|I_j| - 1)\alpha,$$

and $\sum_{j \in J} (|I_j| - 1)\alpha = (n - m)\alpha$ since $|I| = n$ and $|J| = m$. The statement of the theorem now follows since problem (7) is optimized by $\tau^* = \frac{n(n-m)}{(n-1)}$, $\alpha^* = \frac{n}{n-1}$ and $\beta^* = \frac{-1}{n-1}$. \square

The proof of Theorem 2 shows that the upper bound $\sqrt{\frac{n-m}{n-1}}$ on the worst-case type-2 Wasserstein distance $\bar{\mathcal{C}}_2(n, m)$ is tight whenever there is an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ whose scenarios $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n$ correspond to a Gram matrix $\mathbf{S} = [\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n]^\top [\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n] = \frac{n}{n-1} \mathbb{I} - \frac{1}{n-1} \mathbf{e}\mathbf{e}^\top$, which implies $\|\boldsymbol{\xi}_i\|_2 = \sqrt{s_{ii}} = 1$ for all $i \in I$. We now show that such an empirical distribution exists when $d \geq n - 1$.

Proposition 1 For $d \geq n - 1$, there is $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ with $\|\xi\|_2 \leq 1$ for all $\xi \in \text{supp}(\hat{\mathbb{P}}_n)$ such that $C_2(\hat{\mathbb{P}}_n, m) = \sqrt{\frac{n-m}{n-1}}$.

Proof Assume first that $d = n$ and consider the empirical distribution $\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i}$ defined through

$$\xi_i = y\mathbf{e} + (x - y)\mathbf{e}_i \in \mathbb{R}^n \quad \text{with} \quad x = \sqrt{\frac{n-1}{n}} \quad \text{and} \quad y = \frac{-1}{\sqrt{n(n-1)}}. \quad (8)$$

A direct calculation reveals that $\mathbf{S} = [\xi_1, \dots, \xi_n]^\top [\xi_1, \dots, \xi_n] = \frac{n}{n-1} \mathbb{I} - \frac{1}{n-1} \mathbf{e}\mathbf{e}^\top$.

To prove the statement for $d = n - 1$, we note that the n scenarios in (8) lie on the $(n - 1)$ -dimensional subspace \mathcal{H} orthogonal to $\mathbf{e} \in \mathbb{R}^n$. Thus, there exists a rotation that maps \mathcal{H} to $\mathbb{R}^{n-1} \times \{0\}$. As the Gram matrix is invariant under rotations, the rotated scenarios give rise to an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^{n-1}, n)$ satisfying the statement of the proposition. Likewise, for $d > n$ the linear transformation $\xi_i \mapsto (\mathbb{I}, \mathbf{0})^\top \xi_i$, $\mathbb{I} \in \mathbb{R}^{n \times n}$ and $\mathbf{0} \in \mathbb{R}^{n \times (d-n)}$, generates an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ that satisfies the statement of the proposition. \square

Proposition 1 requires that $d \geq n - 1$, which appears to be restrictive. We note, however, that this condition is only sufficient (and not necessary) to guarantee the tightness of the bound from Theorem 2. Moreover, we will observe in Section 2.3 that the bound of Theorem 2 provides surprisingly accurate guidance for the Wasserstein distance between practice-relevant empirical distributions $\hat{\mathbb{P}}_n$ and their continuously reduced optimal distributions.

2.2 Fundamental Limits for the Type-1 Wasserstein Distance

In analogy to the previous section, we now derive a revised upper bound on $\bar{C}_l(n, m)$ for the type-1 Wasserstein distance.

Theorem 3 The worst-case type-1 Wasserstein distance satisfies $\bar{C}_1(n, m) \leq \sqrt{\frac{n-m}{n-1}}$.

Note that this bound is identical to the bound of Theorem 2 for $l = 2$.

Proof of Theorem 3 Leveraging again Theorem 1 and Remark 1, we obtain that

$$\begin{aligned} \bar{C}_1(n, m) &= \max_{\{\xi_i\} \subseteq \mathbb{R}^d} \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{gmed}(I_j)\|_2 \\ \text{s.t.} \quad &\|\xi_i\|_2 \leq 1 \quad \forall i \in I. \end{aligned}$$

We show that $\bar{C}_1(n, m) \leq \bar{C}_2(n, m)$ for all n and $m = 1, \dots, n$, which in turn proves the statement of the theorem by virtue of Theorem 2. To this end, we observe that

$$\begin{aligned} \bar{C}_1(n, m) &\leq \max_{\{\xi_i\} \subseteq \mathbb{R}^d} \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2 \\ \text{s.t.} \quad &\|\xi_i\|_2 \leq 1 \quad \forall i \in I \end{aligned} \quad (9)$$

$$\begin{aligned} &\leq \max_{\{\xi_i\} \subseteq \mathbb{R}^d} \min_{\{I_j\} \in \mathfrak{P}(I, m)} \left[\frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2^2 \right]^{1/2} \\ \text{s.t.} \quad &\|\xi_i\|_2 \leq 1 \quad \forall i \in I, \end{aligned} \quad (10)$$

where the first inequality follows from the definition of the geometric median, which ensures that

$$\sum_{i \in I_j} \|\xi_i - \text{gmed}(I_j)\|_2 \leq \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2 \quad \forall j \in J,$$

and the second inequality is due to the arithmetic-mean quadratic-mean inequality (Steele 2004, Exercise 2.14). The statement of the theorem now follows from the observation that the optimal value of (10) is identical to $\overline{C}_2(n, m)$. \square

In the next proposition we derive a lower bound on $\overline{C}_1(n, m)$.

Proposition 2 *For $d \geq n - 1$, the worst-case type-1 Wasserstein distance satisfies $\overline{C}_1(n, m) \geq \sqrt{\frac{(n-m)(n-m+1)}{n(n-1)}}$.*

Proof Assume first that $d = n$, and consider the empirical distribution $\hat{\mathbb{P}}_n$ with scenarios defined as in (8). Let $\{I_j\} \in \mathfrak{P}(I, m)$ be an arbitrary m -set partition of I and note that $\text{gmed}(I_j) = \text{mean}(I_j)$ for every $j \in J$ due to the permutation symmetry of the ξ_i . This is indeed the case because $\mathbf{0} \in \partial f_j(\text{mean}(I_j))$ for each $f_j(\zeta) = \sum_{i \in I_j} \|y\mathbf{e} + (x - y)\mathbf{e}_i - \zeta\|_2$, $j \in J$. Thus, we have

$$\begin{aligned} & \|\xi_i - \text{mean}(I_j)\|_2 \\ &= \left[\left(x - \frac{x + (|I_j| - 1)y}{|I_j|} \right)^2 + (|I_j| - 1) \left(y - \frac{x + (|I_j| - 1)y}{|I_j|} \right)^2 \right]^{1/2} \\ &= \left[\frac{|I_j| - 1}{|I_j|} \right]^{1/2} (x - y) = \left[\frac{n(|I_j| - 1)}{(n - 1)|I_j|} \right]^{1/2} \quad \forall i \in I_j, \end{aligned}$$

where the last equality follows from the definitions of x and y in (8). By Theorem 1 and Remark 1 we therefore obtain

$$\begin{aligned} C_1(\hat{\mathbb{P}}_n, m) &= \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2 \\ &= \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{\sqrt{n(n-1)}} \sum_{j \in J} \sqrt{|I_j|(|I_j| - 1)}. \end{aligned}$$

By introducing auxiliary variables $z_j = |I_j| - 1 \in \mathbb{N}_0$, $j \in J$, we find that determining $C_1(\hat{\mathbb{P}}_n, m)$ is tantamount to solving

$$C_1(\hat{\mathbb{P}}_n, m) = \frac{1}{\sqrt{n(n-1)}} \min_{\{z_j\} \subseteq \mathbb{N}_0} \left\{ \sum_{j \in J} \sqrt{z_j(z_j + 1)} : \sum_{j \in J} z_j = n - m \right\}.$$

Observe that the objective function of $z_1 = n - m$ and $z_2 = \dots = z_m = 0$ evaluates to $\sqrt{(n - m)(n - m + 1)}$, which implies that $C_1(\hat{\mathbb{P}}_n, m) \leq \sqrt{\frac{(n-m)(n-m+1)}{n(n-1)}}$. Hence,

it remains to establish the reverse inequality. To this end, we note that

$$\begin{aligned}
\sum_{j \in J} \sqrt{z_j(z_j + 1)} &= \left[\sum_{j \in J} z_j(z_j + 1) + \sum_{\substack{j, j' \in J \\ j \neq j'}} \sqrt{z_j z_{j'}(1 + z_j)(1 + z_{j'})} \right]^{1/2} \\
&\geq \left[\sum_{j \in J} z_j(z_j + 1) + \sum_{\substack{j, j' \in J \\ j \neq j'}} z_j z_{j'} \right]^{1/2} \\
&= \left[\left(\sum_{j \in J} z_j \right)^2 + \sum_{j \in J} z_j \right]^{1/2} = \sqrt{(n - m)(n - m + 1)},
\end{aligned}$$

and thus the claim follows for $d = n$. The cases $d = n - 1$ and $d > n$ can be reduced to the case $d = n$ as in Proposition 1. Details are omitted for brevity. \square

Proposition 2 asserts that $\bar{C}_1(n, m) \gtrsim \frac{n-m}{n-1} = \bar{C}_2^2(n, m)$ whenever $d \geq n - 1$. Together with Theorem 3, we thus obtain the following relation between the worst-case Wasserstein distances of types $l = 1$ and $l = 2$:

$$\bar{C}_2^2(n, m) \leq \bar{C}_1(n, m) \leq \bar{C}_2(n, m).$$

We conjecture that the lower bound is tighter, but we were not able to prove this.

2.3 Discussion

Theorems 2 and 3 imply that $\bar{C}_l(n, m) \lesssim \sqrt{1 - p}$ for large n and for $l \in \{1, 2\}$, where $p = \frac{m}{n}$ represents the desired reduction factor. The significance of this result is that it offers *a priori* guidelines for selecting the number m of support points in the reduced distribution. To see this, consider any empirical distribution $\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i \in I} \delta_{\xi_i}$, and denote by $r \geq 0$ and $\boldsymbol{\mu} \in \mathbb{R}^d$ the radius and the center of any (ideally the smallest) ball enclosing ξ_1, \dots, ξ_n , respectively. In this case, we have

$$C_l(\hat{\mathbb{P}}_n, m) = r \cdot C_l\left(\frac{1}{n} \sum_{i \in I} \delta_{\frac{\xi_i - \boldsymbol{\mu}}{r}}, m\right) \leq r \cdot \bar{C}_l(n, m) \lesssim r \cdot \sqrt{1 - p}, \quad (11)$$

where the inequality holds because $\|(\xi_i - \boldsymbol{\mu})/r\|_2 \leq 1$ for every $i \in I$. Note that (11) enables us to find an upper bound on the smallest m guaranteeing that $C_l(\hat{\mathbb{P}}_n, m)$ falls below a prescribed threshold (*i.e.*, guaranteeing that the reduced m -point distribution remains within some prescribed distance from $\hat{\mathbb{P}}_n$).

Even though the inequality in (11) can be tight, which has been established in Proposition 1, one might suspect that typically $C_l(\hat{\mathbb{P}}_n, m)$ is significantly smaller than $r \cdot \sqrt{1 - p}$ when the points $\xi_i \in \mathbb{R}^d$, $i \in I$, are sampled randomly from a standard distribution, *e.g.*, a multivariate uniform or normal distribution. However, while the upper bound (11) can be loose for low-dimensional data, Proposition 3 below suggests that it is surprisingly tight in high dimensions—at least for $l = 2$.

Proposition 3 For any $\epsilon > 0$ and $\delta > 0$ there exist $c > 0$ and $d \in \mathbb{N}$ such that

$$\mathbb{P}^n \left(\|\xi_i\|_2 \leq 1 \ \forall i \in I \text{ and } C_2 \left(\frac{1}{n} \sum_{i \in I} \delta_{\xi_i}, m \right) \geq \sqrt{1-p} - \delta \right) \geq 1 - \epsilon, \quad (12)$$

where $p = \frac{m}{n}$, and the support points ξ_i , $i \in I$, are sampled independently from the normal distribution \mathbb{P} with mean $\mathbf{0} \in \mathbb{R}^d$ and covariance matrix $(\sqrt{d-1} + c)^{-2} \mathbb{I} \in \mathbb{S}^d$.

Proposition 3 can be paraphrased as follows. Sampling the ξ_i , $i \in I$, independently from the normal distribution \mathbb{P} yields a (random) empirical distribution $\hat{\mathbb{P}}_n$ that is feasible and δ -suboptimal in (1) with probability $1 - \epsilon$. The intuition behind this result is that, in high dimensions, samples drawn from \mathbb{P} are almost orthogonal and close to the surface of the unit ball with high probability. Indeed, these two properties are shared by the worst case distribution (8) in high dimensions.

Proof of Proposition 3 Theorem 1 and Remark 1 imply that

$$C_2^2 \left(\frac{1}{n} \sum_{i \in I} \delta_{\xi_i}, m \right) = \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2^2. \quad (13)$$

From the proof of Theorem 2 we further know that (13) can be expressed as a continuous function $f(\mathbf{S})$ of the Gram matrix $\mathbf{S} = [\xi_1, \dots, \xi_n]^\top [\xi_1, \dots, \xi_n]$, that is,

$$f(\mathbf{S}) = \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{n} \sum_{j \in J} \frac{1}{|I_j|^2} \sum_{i \in I_j} \left(|I_j|^2 s_{ii} - 2|I_j| \sum_{k \in I_j} s_{ik} + \sum_{k \in I_j} s_{kk} + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} s_{kk'} \right).$$

An elementary calculation shows that $f(\mathbb{I}) = \frac{n-m}{n} = 1-p$. Thus, by the continuity of $f(\cdot)$, there exists $\eta \in (0, 1)$ with $\sqrt{f(\mathbf{S})} \geq \sqrt{1-p} - \delta$ whenever $\|\mathbf{S} - \mathbb{I}\|_{\max} \leq \eta$.

We are now ready to construct \mathbb{P} . First, select $c > 0$ large enough to ensure that

$$\left(1 - \frac{4}{c^2} e^{-\frac{c^2}{4}} \right)^n \geq 1 - \frac{\epsilon}{2}.$$

Then, select $d \in \mathbb{N}$ large enough such that

$$\frac{\sqrt{d-1} - c}{\sqrt{d-1} + c} \geq 1 - \eta \quad \text{and} \quad n(n-1)\Phi(-\eta(\sqrt{d-1} + c)) \leq \frac{\epsilon}{2},$$

where $\Phi(\cdot)$ denotes the univariate standard normal distribution function. Observe that the distribution \mathbb{P} is completely determined by c and d . Next, we find that

$$\begin{aligned} & \mathbb{P}^n \left(\|\xi_i\|_2 \leq 1 \ \forall i \text{ and } C_2 \left(\frac{1}{n} \sum_{i \in I} \delta_{\xi_i}, m \right) \geq \sqrt{1-p} - \delta \right) \\ & \geq \mathbb{P}^n (\|\xi_i\|_2 \leq 1 \ \forall i \text{ and } \|\mathbf{S} - \mathbb{I}\|_{\max} \leq \eta) \\ & = \mathbb{P}^n \left(1 - \eta \leq \|\xi_i\|_2 \leq 1 \ \forall i \text{ and } |\xi_i^\top \xi_j| \leq \eta \ \forall i \neq j \right) \\ & \geq \mathbb{P}^n \left(\frac{\sqrt{d-1} - c}{\sqrt{d-1} + c} \leq \|\xi_i\|_2 \leq 1 \ \forall i \text{ and } |\xi_i^\top \xi_j| \leq \eta \cdot \|\xi_j\|_2 \ \forall i \neq j \right) \\ & \geq \mathbb{P}^n \left(\frac{\sqrt{d-1} - c}{\sqrt{d-1} + c} \leq \|\xi_i\|_2 \leq 1 \ \forall i \right) + \mathbb{P}^n \left(|\xi_i^\top \xi_j| \leq \eta \cdot \|\xi_j\|_2 \ \forall i \neq j \right) - 1, \quad (14) \end{aligned}$$

where the first and second inequalities follow from the construction of η and c , respectively, while the third inequality exploits the union bound. We also have

$$\begin{aligned} \mathbb{P}^n \left(\frac{\sqrt{d-1}-c}{\sqrt{d-1}+c} \leq \|\xi_i\|_2 \leq 1 \ \forall i \right) &= \mathbb{P} \left(\frac{\sqrt{d-1}-c}{\sqrt{d-1}+c} \leq \|\xi_1\|_2 \leq 1 \right)^n \\ &\geq \left(1 - \frac{4}{c^2} e^{-\frac{c^2}{4}} \right)^n \geq 1 - \frac{\epsilon}{2}, \end{aligned} \quad (15)$$

where the equality holds due to the independence of the ξ_i , while the first and second inequalities follow from Lemma 2.8 in Hopcroft and Kannan (2012) and the construction of c , respectively. By the choice of d , we finally obtain

$$\begin{aligned} \mathbb{P}^n \left(|\xi_i^\top \xi_j| \leq \eta \cdot \|\xi_j\|_2 \ \forall i \neq j \right) &\geq 1 - \sum_{i \neq j} \mathbb{P}^n \left(|\xi_i^\top \xi_j| \geq \eta \cdot \|\xi_j\|_2 \right) \\ &= 1 - n(n-1) \Phi(-\eta(\sqrt{d-1}+c)) \geq 1 - \frac{\epsilon}{2}. \end{aligned} \quad (16)$$

The equality in (16) holds due to the rotation symmetry of \mathbb{P} , which implies that

$$\mathbb{P}^n \left(|\xi_i^\top \xi_j| \geq \eta \cdot \|\xi_j\|_2 \right) = \mathbb{P} \left(|\xi_1^\top \mathbf{e}_1| \geq \eta \right) = 2\Phi(-\eta(\sqrt{d-1}+c)).$$

The claim then follows by substituting (15) and (16) into (14). \square

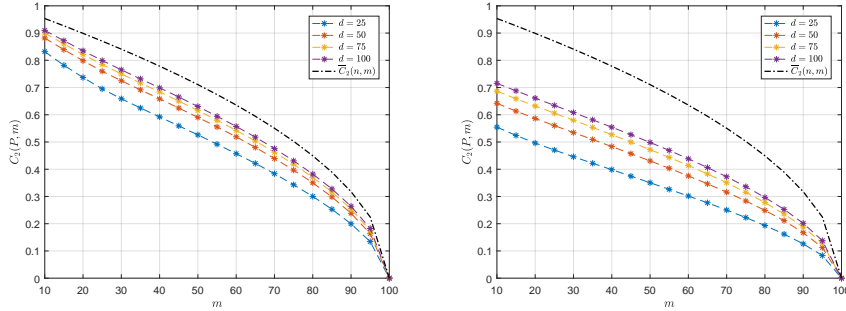


Fig. 1: Comparison between $C_2(\hat{\mathbb{P}}_n, m)$ and $\overline{C}_2(n, m)$ under uniform (left panel) and normal (right panel) sampling.

Figure 1 compares $\overline{C}_2(m, n)$ with $C_2(\hat{\mathbb{P}}_n, m)$ for $n = 100$, $m \in \{10, \dots, 100\}$ and $d \in \{25, 50, 75, 100\}$. The n original support points are sampled randomly from the uniform distribution on the unit ball (left panel) and the normal distribution from Proposition 3 with $c = 2.97$, which ensures that $\|\xi_i\|_2 \leq 1$ with 95% probability (right panel). Note that $C_2(\mathbb{P}_n, m)$ is random. Thus, all shown values are averaged across 100 independent trials. Figure 1 confirms that $C_2(\hat{\mathbb{P}}_n, m)$ approaches the worst-case bound $\overline{C}_2(n, m)$ as the dimension d increases.

3 Guarantees for Discrete Scenario Reduction

For n -point empirical distributions $\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i}$ supported on \mathbb{R}^d , we now study the loss of optimality incurred by solving the discrete scenario reduction problem instead of its continuous counterpart. More precisely, we want to determine the point-wise largest lower bound $\underline{\kappa}_l(n, m)$ and the point-wise smallest upper bound $\bar{\kappa}_l(n, m)$ that satisfy

$$\underline{\kappa}_l(n, m) \cdot C_l(\hat{\mathbb{P}}_n, m) \leq D_l(\hat{\mathbb{P}}_n, m) \leq \bar{\kappa}_l(n, m) \cdot C_l(\hat{\mathbb{P}}_n, m) \quad \forall \hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n) \quad (17)$$

for the Wasserstein distances of type $l \in \{1, 2\}$. Note that the existence of finite bounds $\underline{\kappa}_l(n, m)$ and $\bar{\kappa}_l(n, m)$ is not a priori obvious as they do not depend on the dimension d . Moreover, while it is clear that $\underline{\kappa}_l(n, m) \geq 1$ if it exists, it does not seem easy to derive a naïve upper bound on $\bar{\kappa}_l(n, m)$.

Our derivations in this section will use the following result, which is the analogue of Theorem 1 for the discrete scenario reduction problem.

Theorem 4 *For any type- l Wasserstein distance induced by any norm $\|\cdot\|$, the discrete scenario reduction problem can be reformulated as*

$$D_l(\hat{\mathbb{P}}_n, m) = \min_{\{I_j\} \in \mathfrak{P}(I, m)} \left[\frac{1}{n} \sum_{j \in J} \min_{\zeta_j \in \{\xi_i : i \in I_j\}} \sum_{i \in I_j} \|\xi_i - \zeta_j\|^l \right]^{1/l}.$$

Proof The proof is similar to the proof of Theorem 1 and is therefore omitted. \square

The remainder of this section derives lower and upper bounds on $\underline{\kappa}_l(n, m)$ and $\bar{\kappa}_l(n, m)$ for Wasserstein distances of type $l = 2$ (Section 3.1) and $l = 1$ (Section 3.2), respectively. To eliminate trivial cases, we assume throughout this section that $n \geq 2$, $m \in \{1, \dots, n-1\}$ and $d \geq 2$.

3.1 Guarantees for the Type-2 Wasserstein Distance

We first bound $\bar{\kappa}_2(n, m)$ in equation (17) from above (Theorem 5) and below (Proposition 4).

Theorem 5 *The upper bound $\bar{\kappa}_2(n, m)$ in (17) satisfies $\bar{\kappa}_2(n, m) \leq \sqrt{2}$ for all n, m .*

Proof The proof proceeds in two steps. We first show that $\bar{\kappa}_2(n, m) \leq \sqrt{2}$ for all n when $m = 1$ (Step 1). Then we extend the result to all n and m (Step 2).

Step 1: Fix any $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$. W.l.o.g., we can assume that $\text{mean}(I) = \mathbf{0}$ and $\frac{1}{n} \sum_{i \in I} \|\xi_i\|_2^2 = 1$ by re-positioning and scaling the atoms ξ_i appropriately. Note that the re-positioning does not affect $C_2(\hat{\mathbb{P}}_n, 1)$ or $D_2(\hat{\mathbb{P}}_n, 1)$, and the positive homogeneity of the Wasserstein distance implies that the scaling affects both $C_2(\hat{\mathbb{P}}_n, 1)$ and $D_2(\hat{\mathbb{P}}_n, 1)$ in the same way and thus preserves their ratio $\bar{\kappa}_2(n, 1)$. Theorem 1 and Remark 1 then imply that

$$C_2(\hat{\mathbb{P}}_n, 1) = \left[\frac{1}{n} \sum_{i \in I} \|\xi_i - \text{mean}(I)\|_2^2 \right]^{1/2} = 1.$$

Step 1 is thus complete if we can show that $D_2(\hat{\mathbb{P}}_n, 1) \leq \sqrt{2}$. Indeed, we have

$$\begin{aligned}
D_2^2(\hat{\mathbb{P}}_n, 1) &= \min_{j \in I} \frac{1}{n} \sum_{i \in I} \|\xi_i - \xi_j\|_2^2 = \min_{j \in I} \frac{1}{n} \sum_{i \in I} (\xi_i - \xi_j)^\top (\xi_i - \xi_j) \\
&= \min_{j \in I} \frac{1}{n} \sum_{i \in I} (\xi_i^\top \xi_i - 2\xi_i^\top \xi_j + \xi_j^\top \xi_j) = \min_{j \in I} \frac{1}{n} \sum_{i \in I} (\xi_i^\top \xi_i + \xi_j^\top \xi_j) \\
&= \min_{j \in I} \|\xi_j\|_2^2 + \frac{1}{n} \sum_{i \in I} \|\xi_i\|_2^2 = 1 + \min_{j \in I} \|\xi_j\|_2^2 \leq 2,
\end{aligned} \tag{18}$$

where the first equality is due to Theorem 4, the fourth follows from $\sum_{i \in I} \xi_i = n \cdot \text{mean}(I) = \mathbf{0}$, and the inequality holds since $\min_{j \in I} \|\xi_j\|_2^2 \leq \frac{1}{n} \sum_{i \in I} \|\xi_i\|_2^2 = 1$.

Step 2: Fix any $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$. Theorem 1 and Remark 1 imply that

$$C_2(\hat{\mathbb{P}}_n, m) = \min_{\{I_j\} \in \mathfrak{P}(I, m)} \left[\frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{mean}(I_j)\|_2^2 \right]^{1/2}.$$

For an optimal partition $\{I_j^*\}$ to this problem, $C_2(\hat{\mathbb{P}}_n, m)$ can be expressed as

$$C_2(\hat{\mathbb{P}}_n, m) = \left[\sum_{j \in J} \frac{|I_j^*|}{n} C_{2,j}^2 \right]^{1/2} \quad \text{with} \quad C_{2,j} = \left[\frac{1}{|I_j^*|} \sum_{i \in I_j^*} \|\xi_i - \text{mean}(I_j^*)\|_2^2 \right]^{1/2}.$$

From our discussion in Step 1 we know that $C_{2,j}$ represents the type-2 Wasserstein distance between the conditional empirical distribution $\hat{\mathbb{P}}_n^j = \frac{1}{|I_j^*|} \sum_{i \in I_j^*} \delta_{\xi_i}$ and its closest Dirac distribution, that is, $C_2(\hat{\mathbb{P}}_n^j, 1)$. Analogously, we obtain that

$$\begin{aligned}
D_2(\hat{\mathbb{P}}_n, m) &\leq \left[\sum_{j \in J} \frac{|I_j^*|}{n} D_{2,j}^2 \right]^{1/2} \quad \text{with} \quad D_{2,j} = \left[\min_{j \in I_j^*} \frac{1}{|I_j^*|} \sum_{i \in I_j^*} \|\xi_i - \xi_j\|_2^2 \right]^{1/2} \\
&\leq \left[\sum_{j \in J} \frac{|I_j^*|}{n} (2C_{2,j}^2) \right]^{1/2} = \sqrt{2} C_2(\hat{\mathbb{P}}_n, m),
\end{aligned}$$

where the first inequality holds since the optimal partition $\{I_j^*\}$ for $C_2(\hat{\mathbb{P}}_n, m)$ is typically suboptimal in $D_2(\hat{\mathbb{P}}_n, m)$, the second inequality follows from the fact that $D_{2,j} = D_2(\hat{\mathbb{P}}_n^j, 1)$ and $D_2(\hat{\mathbb{P}}_n^j, 1) \leq \sqrt{2} C_2(\hat{\mathbb{P}}_n^j, 1)$ due to Step 1, and the identity follows from the definition of $C_{2,j}$. The statement now follows. \square

Proposition 4 *There is $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ with $D_2(\hat{\mathbb{P}}_n, m) = \sqrt{2} C_2(\hat{\mathbb{P}}_n, m)$ for all n, m .*

Proof In analogy to the proof of Theorem 5, we first show the statement for $m = 1$ (Step 1) and then extend the result to $m > 1$ (Step 2).

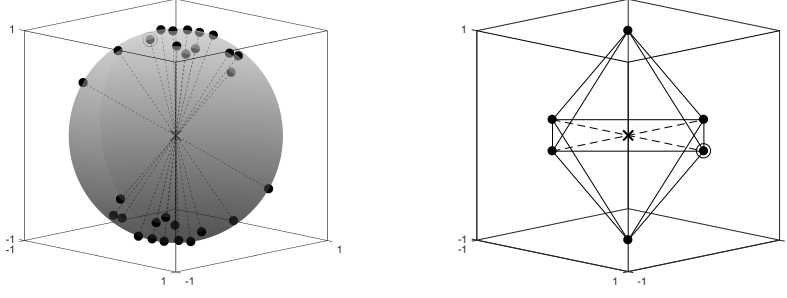


Fig. 2: Empirical distributions in \mathbb{R}^3 that maximize the ratio between $D_l(\hat{\mathbb{P}}_n, 1)$ and $C_l(\hat{\mathbb{P}}_n, 1)$ for $l = 2$ and $\|\cdot\| = \|\cdot\|_2$ (left panel) as well as $l = 1$ and $\|\cdot\| = \|\cdot\|_1$ (right panel). In both cases, the continuous scenario reduction problem is optimized by the Dirac distribution at $\mathbf{0}$ (marked as \times), whereas the discrete scenario reduction problem is optimized by any of the atoms (such as \circ).

Step 1: The first step in the proof of Theorem 5 shows that $\hat{\mathbb{P}}_n \in \mathcal{P}_{\mathbb{E}}(\mathbb{R}^d, n)$ satisfies $D_2(\hat{\mathbb{P}}_n, 1) = \sqrt{2}C_2(\hat{\mathbb{P}}_n, 1)$ if $\sum_{i \in I} \xi_i = \mathbf{0}$ and $\|\xi_1\| = \dots = \|\xi_n\| = 1$. For an even number $n = 2k$, $k \in \mathbb{N}$, both conditions are satisfied if we place ξ_1, \dots, ξ_k on the surface of the unit ball in \mathbb{R}^d and then choose $\xi_{k+i} = -\xi_i$ for $i = 1, \dots, k$ (see left panel of Figure 2 for an illustration in \mathbb{R}^3). Likewise, for an odd number $n = 2k+3$, $k \in \mathbb{N}_0$, we can place ξ_1, \dots, ξ_k on the surface of the unit ball, choose $\xi_{k+i} = -\xi_i$ for $i = 1, \dots, k$ and fix $\xi_{2k+1} = \mathbf{e}_1$, $\xi_{2k+2} = -\frac{1}{2}\mathbf{e}_1 + \frac{\sqrt{3}}{2}\mathbf{e}_2$ and $\xi_{2k+3} = -\frac{1}{2}\mathbf{e}_1 - \frac{\sqrt{3}}{2}\mathbf{e}_2$.

Step 2: To prove the statement for $m > 1$, we construct an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_{\mathbb{E}}(\mathbb{R}^d, n)$ whose atoms satisfy $\text{supp}(\hat{\mathbb{P}}_n) = \Xi_1 \cup \Xi_2$ with $|\Xi_1| = n - m + 1$ and $|\Xi_2| = m - 1$. The atoms $\xi_1, \dots, \xi_{n-m+1}$ in Ξ_1 are selected according to the recipe outlined in Step 1, whereas the atoms $\xi_{n-m+2}, \dots, \xi_n$ in Ξ_2 satisfy $\xi_{n-m+1+i} = (1 + iM)\mathbf{e}_1$, $i = 1, \dots, m - 1$, for any number M satisfying $M > 2\sqrt{n - m + 1}$. A direct calculation then shows that the atoms in Ξ_2 are sufficiently far away from those in Ξ_1 as well as from each other so that any optimal partition $\{I_j^*\}$ to the discrete scenario reduction problem in Theorem 4 as well as the continuous scenario reduction problem in Theorem 1 consists of the sets $\{i : \xi_i \in \Xi_1\}$ and $\{i\}$, $\xi_i \in \Xi_2$. The result then follows from the fact that either problem accumulates a Wasserstein distance of 0 over the atoms in Ξ_2 , whereas the Wasserstein distance of $D_2(\hat{\mathbb{P}}_n, m)$ is a factor of $\sqrt{2}$ bigger than the Wasserstein distance of $C_2(\hat{\mathbb{P}}_n, m)$ over the atoms in Ξ_1 (see Step 1). \square

Theorem 5 and Proposition 4 imply that $\bar{\kappa}_2(n, m) = \sqrt{2}$ for all n and m , that is, the bound is indeed *independent* of both the number of atoms n in the empirical distribution and the number of atoms m in the reduced distribution. We now show that the naïve lower bound of 1 on the approximation ratio is essentially tight.

Proposition 5 *The lower bound $\underline{\kappa}_2(n, m)$ in (17) satisfies $\underline{\kappa}_2(n, m) = 1$ whenever $n \geq 3$ and $m \in \{1, \dots, n - 2\}$, while $\underline{\kappa}_2(n, n - 1) = \sqrt{2}$ always.*

Proof We first prove $\underline{\kappa}_2(n, m) = 1$ when $m = 1$ and $n \geq 3$ (Step 1) and when $m \in \{2, \dots, n - 2\}$ (Step 2). Then, we show $\underline{\kappa}_2(n, n - 1) = \sqrt{2}$ (Step 3).

Step 1: Choose $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ such that the first $n-1$ atoms ξ_1, \dots, ξ_{n-1} are selected according to the recipe outlined in Step 1 in the proof of Proposition 4 and $\xi_n = \mathbf{0}$. We thus have $\text{mean}(I) = \mathbf{0}$, and Theorem 1 and Remark 1 imply that the optimal continuous scenario reduction is given by the Dirac distribution $\delta_{\mathbf{0}}$. Since $\mathbf{0} \in \text{supp}(\hat{\mathbb{P}}_n)$, we have $C_2(\hat{\mathbb{P}}_n, 1) = D_2(\hat{\mathbb{P}}_n, 1)$ and the result follows.

Step 2: To prove the statement for $m > 1$, we proceed as in Step 2 in the proof of Proposition 4. In particular, we construct an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ whose atoms satisfy $\text{supp}(\hat{\mathbb{P}}_n) = \Xi_1 \cup \Xi_2$ with $|\Xi_1| = n - m + 1$ and $|\Xi_2| = m - 1$. The atoms $\xi_1, \dots, \xi_{n-m+1}$ in Ξ_1 are selected according to the recipe outlined in Step 1 of this proof, whereas the remaining atoms $\xi_{n-m+1}, \dots, \xi_n$ in Ξ_2 satisfy $\xi_{n-m+1+i} = (1 + iM)\mathbf{e}_1$, $i = 1, \dots, m - 1$, for any $M > 2\sqrt{n - m + 1}$. A similar argument as in the proof of Proposition 4 then shows that $C_2(\hat{\mathbb{P}}_n, m) = D_2(\hat{\mathbb{P}}_n, m)$.

Step 3: Fix any $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$. W.l.o.g., assume that $\{\xi_{n-1}, \xi_n\}$ is the closest pair of atoms in terms of Euclidean distance, and let $d_{\min} = \|\xi_n - \xi_{n-1}\|_2$. One readily verifies that the partition $I_j^* = \{j\}$, $j = 1, \dots, n - 2$, and $I_{n-1}^* = \{n - 1, n\}$ optimizes both the discrete scenario reduction problem in Theorem 4 as well as the continuous scenario reduction problem in Theorem 1. We thus have $C_2(\hat{\mathbb{P}}_n, n - 1) = \frac{1}{\sqrt{2n}}d_{\min}$ and $D_2(\hat{\mathbb{P}}_n, n - 1) = \frac{1}{\sqrt{n}}d_{\min}$, which concludes the proof. \square

Hence, for any empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ the type-2 Wasserstein distance between the minimizer of the *discrete* scenario reduction problem and $\hat{\mathbb{P}}_n$ exceeds the Wasserstein distance between the minimizer of the *continuous* scenario reduction problem and $\hat{\mathbb{P}}_n$ by up to 41.4%, and the bound is attainable for any n, m .

3.2 Guarantees for the Type-1 Wasserstein Distance

In analogy to Section 3.1, we first bound $\bar{\kappa}_1(n, m)$ from above (Theorem 6) and below (Proposition 6). In contrast to the previous section, we consider an arbitrary norm $\|\cdot\|$, and we adapt the definition of the geometric median accordingly.

Theorem 6 *The upper bound $\bar{\kappa}_1(n, m)$ in (17) satisfies $\bar{\kappa}_1(n, m) \leq 2$ whenever $m \in \{2, \dots, n - 2\}$ as well as $\bar{\kappa}_1(n, 1) \leq 2(1 - \frac{1}{n})$ and $\bar{\kappa}_1(n, n - 1) \leq 1$.*

Proof We first prove the statement for $m = 1$ (Step 1) and then extend the result to $m \in \{2, \dots, n - 2\}$ (Step 2) and $m = n - 1$ (Step 3).

Step 1: Fix any $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$. As in the proof of Theorem 5, we can assume that $\text{gmed}(I) = \mathbf{0}$ and $\frac{1}{n} \sum_{i \in I} \|\xi_i\| = 1$ by re-positioning and scaling the atoms ξ_i appropriately. Theorem 1 and Remark 1 then imply that for $m = 1$, we have

$$C_1(\hat{\mathbb{P}}_n, 1) = \frac{1}{n} \sum_{i \in I} \|\xi_i - \text{gmed}(I)\| = 1.$$

Step 1 is thus complete if we can show that $D_1(\hat{\mathbb{P}}_n, 1) \leq 2(1 - \frac{1}{n})$. Indeed, we have

$$\begin{aligned} D_1(\hat{\mathbb{P}}_n, 1) &= \min_{j \in I} \frac{1}{n} \sum_{i \in I} \|\xi_i - \xi_j\| = \min_{j \in I} \frac{1}{n} \sum_{i \in I \setminus \{j\}} \|\xi_i - \xi_j\| \\ &\leq \min_{j \in I} \frac{1}{n} \sum_{i \in I \setminus \{j\}} (\|\xi_i\| + \|\xi_j\|) = \min_{j \in I} \frac{1}{n} \left((n-2)\|\xi_j\| + \sum_{i \in I} \|\xi_i\| \right) \\ &= \min_{j \in I} \frac{1}{n} ((n-2)\|\xi_j\| + n) = 1 + \frac{n-2}{n} \cdot \min_{j \in I} \|\xi_j\| \leq 2 \left(1 - \frac{1}{n} \right), \end{aligned}$$

where the two inequalities follow from the triangle inequality and the fact that $\min_{j \in I} \|\xi_j\| \leq \frac{1}{n} \sum_{i \in I} \|\xi_i\| = 1$, respectively.

Step 2: Fix any $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$. Theorem 1 and Remark 1 then imply that

$$C_1(\hat{\mathbb{P}}_n, m) = \min_{\{I_j\} \in \mathfrak{P}(I, m)} \frac{1}{n} \sum_{j \in J} \sum_{i \in I_j} \|\xi_i - \text{gmed}(I_j)\|.$$

Let $\{I_j^*\}$ be an optimal partition for this problem. The same arguments as in the proof of Theorem 5 show that

$$\begin{aligned} D_1(\hat{\mathbb{P}}_n, m) &\leq \sum_{j \in J} \frac{|I_j^*|}{n} D_{1,j} \quad \text{with } D_{1,j} = \min_{j \in I_j^*} \frac{1}{|I_j^*|} \sum_{i \in I_j^*} \|\xi_i - \xi_j\| \\ &\leq \sum_{j \in J} \frac{|I_j^*|}{n} (2C_{1,j}) \quad \text{with } C_{1,j} = \frac{1}{|I_j^*|} \sum_{i \in I_j^*} \|\xi_i - \text{gmed}(I_j^*)\|, \end{aligned}$$

and the last expression is equal to $2C_1(\hat{\mathbb{P}}_n, m)$ by definition of $C_{1,j}$.

Step 3: For $n = 2$ and $m = n - 1 = 1$, Step 1 shows that $\bar{\kappa}_1(2, 1) \leq 2(1 - \frac{1}{2}) = 1$. For $n > 2$ and $m = n - 1$, the statement can be derived in the same way as the third step in the proof of Proposition 5. We omit the details for the sake of brevity. \square

Proposition 6 *There is $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ such that $D_1(\hat{\mathbb{P}}_n, m) = 2(1 - \frac{m}{n})C_1(\hat{\mathbb{P}}_n, m)$ under the 1-norm for all n divisible by $2m$, all m and all $d \geq \frac{n}{2m}$.*

Proof We first prove the statement for $m = 1$ (Step 1) and then extend the result to $m > 1$ (Step 2). Throughout the proof, we set $k = \frac{n}{2m}$ and consider w.l.o.g. the case where $d = k$.

Step 1: Fix $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ with the atoms $\xi_i = +\mathbf{e}_i$ as well as $\xi_{k+i} = -\mathbf{e}_i$, $i = 1, \dots, k$. The symmetric placement of the atoms implies that $\text{gmed}(I) = \mathbf{0}$ and hence $C_1(\hat{\mathbb{P}}_n, 1) = 1$. Furthermore, we note that $\|\xi_i - \xi_j\|_1 = 2$ for all $i \neq j$, that is, any two atoms are equidistant from another (see right panel of Figure 2 for an illustration in \mathbb{R}^3). By Theorem 4, any 1-point discrete scenario reduction results in a Wasserstein distance of $2\frac{n-1}{n}$ to $\hat{\mathbb{P}}_n$.

Step 2: To prove the statement for $m > 1$, we construct an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ whose atoms satisfy $\text{supp}(\hat{\mathbb{P}}_n) = \bigcup_{j=1}^m (\Xi_j^+ \cup \Xi_j^-)$ with $|\Xi_j^+| = |\Xi_j^-| = k$, $j = 1, \dots, m$. The atoms $\xi_{2(j-1)k+1}, \dots, \xi_{2(j-1)k+k}$ in Ξ_j^+ satisfy $\xi_{2(j-1)k+i} = +\mathbf{e}_i + jM\mathbf{e}_1$, $i = 1, \dots, k$, whereas the atoms $\xi_{2(j-1)k+k+1}, \dots, \xi_{2jk}$ in Ξ_j^- satisfy $\xi_{2(j-1)k+k+i} = -\mathbf{e}_i + jM\mathbf{e}_1$, $i = 1, \dots, k$, for any number M satisfying $M > 2n + 2$. The same arguments as in the proof of Theorem 5 show that any optimal partition $\{I_j^*\}$ to the discrete scenario reduction problem in Theorem 4 as well as the continuous scenario reduction problem in Theorem 1 consists of the sets indexing the atoms in $\Xi_j^+ \cup \Xi_j^-$, $j = 1, \dots, m$. Step 1 shows that the continuous scenario reduction problem accumulates a Wasserstein distance of 1 over each set, whereas the discrete scenario reduction problem accumulates a Wasserstein distance of $2 \frac{2k-1}{2k}$ over each set. The result then follows from the fact that there are m such sets and hence the ratio of the respective overall Wasserstein distances amounts to $m(2 \frac{2k-1}{2k})/m = 2(1 - \frac{m}{n})$. \square

Theorem 6 and Proposition 6 imply that $\bar{\kappa}_1(n, m) \in [2(1 - m/n), 2]$ for all n and $m \in \{2, \dots, n-2\}$. For the small ratios $m : n$ commonly used in practice, we thus conclude that the bound is *essentially independent* of both the number of atoms n in the empirical distribution and the number of atoms m in the reduced distribution. We close with an analysis of the lower bound $\underline{\kappa}_1(n, m)$.

Proposition 7 *The lower bound $\underline{\kappa}_1(n, m)$ in (17) satisfies $\underline{\kappa}_1(n, m) = 1$ for all n, m .*

Proof The proof widely parallels that of Proposition 5, with the difference that the atoms ξ_1, \dots, ξ_n of the empirical distribution $\hat{\mathbb{P}}_n$ are placed such that a geometric median (as opposed to the mean) of each subset in the optimal partition coincides with one of the atoms in that subset. This allows both continuous and discrete scenario reduction to choose the same support points for the reduced distribution, hence incurring the same Wasserstein distance. Details are omitted for brevity. \square

In conclusion, for any $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ the type-1 Wasserstein distance between the minimizer of the *discrete* scenario reduction problem and $\hat{\mathbb{P}}_n$ exceeds the Wasserstein distance between the minimizer of the *continuous* scenario reduction problem and $\hat{\mathbb{P}}_n$ by up to 100%, and this bound is asymptotically attained for decreasing ratios $m : n$.

4 Solution Methods

We now review existing and propose new solution schemes for the discrete and continuous scenario reduction problems. More precisely, we will study two heuristics for discrete and continuous scenario reduction, respectively, that do not come with approximation guarantees (Section 4.1), we will propose a constant-factor approximation scheme for both the discrete and the continuous scenario reduction problem (Section 4.2), and we will discuss two exact reformulations of these problems as mixed-integer optimization problems (Section 4.3).

In the remainder of this section, we denote by $D_l(\hat{\mathbb{P}}_n, \Xi)$ the type- l Wasserstein distance between $\hat{\mathbb{P}}_n$ and its closest distribution supported on the finite set Ξ . Moreover, for an algorithm providing an upper bound $\bar{D}_l(\hat{\mathbb{P}}_n, m)$ on the discrete scenario reduction problem in \mathbb{R}^d , we define the algorithm's *approximation ratio* as

the maximum fraction $\overline{D}_l(\hat{\mathbb{P}}_n, m)/D_l(\hat{\mathbb{P}}_n, m)$, where the maximum is taken over all n and m , as well as all empirical distributions $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$.

4.1 Heuristics for the Discrete Scenario Reduction Problem

We review in Section 4.1.1 a popular heuristic for the discrete scenario reduction problem due to Dupačová et al (2003). We will show that despite the simplicity and efficiency of the algorithm, there is no finite upper bound on the algorithm's approximation ratio. In Section 4.1.2 we adapt a widely used clustering heuristic to the continuous scenario reduction problem, and we show that this algorithm's approximation ratio cannot be bounded from above either.

4.1.1 Dupačová et al.'s Algorithm

We outline Dupačová et al.'s algorithm for the problem $D_l(\hat{\mathbb{P}}_n, m)$ below.

DUPAČOVÁ ET AL.'S ALGORITHM FOR $D_l(\hat{\mathbb{P}}_n, m)$:
1. Initialize the set of atoms in the reduced set as $R \leftarrow \emptyset$.
2. Select the next atom to be added to the reduced set as
$\zeta \in \arg \min_{\zeta \in \text{supp}(\hat{\mathbb{P}}_n)} D_l(\hat{\mathbb{P}}_n, R \cup \{\zeta\})$
and update $R \leftarrow R \cup \{\zeta\}$.
3. Repeat Step 2 until $ R = m$.

Given an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$, the algorithm iteratively populates the reduced set R containing the atoms of the reduced distribution \mathbb{Q} . Each atom $\zeta \in \text{supp}(\hat{\mathbb{P}}_n)$ is selected greedily so as to minimize the Wasserstein distance between $\hat{\mathbb{P}}_n$ and the closest distribution supported on the augmented reduced set $R \cup \{\zeta\}$. After termination, the distribution \mathbb{Q} can be recovered from the reduced set R as follows. Let $\{I_\zeta\} \in \mathfrak{P}(I, m)$ be any partition of $\text{supp}(\hat{\mathbb{P}}_n)$ into sets I_ζ , $\zeta \in R$, such that I_ζ contains all elements of $\text{supp}(\hat{\mathbb{P}}_n)$ that are closest to ζ (ties may be broken arbitrarily). Then $\mathbb{Q} = \sum_{\zeta \in R} q_\zeta \delta_\zeta$, where $q_\zeta = |I_\zeta|/n$.

Theorem 7 *For every $d \geq 2$ and $l, p \geq 1$, the approximation ratio Dupačová et al.'s algorithm is unbounded.*

Proof The proof constructs a specific distribution $\hat{\mathbb{P}}_n$ (Step 1), bounds $D_l(\hat{\mathbb{P}}_n, m)$ from above (Step 2) and bounds the Wasserstein distance between $\hat{\mathbb{P}}_n$ and the output \mathbb{Q} of Dupačová et al.'s algorithm from below (Step 3).

Step 1: Fix $d \geq 2$, $l, p \geq 1$ and $m = 4$, and consider the empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$ with $n = 4z + 1$ for some positive integer z as well as $\text{supp}(\hat{\mathbb{P}}_n) = \Xi_1 \cup \dots \cup \Xi_4 \cup \{\xi_{4z+1}\}$ with $\Xi_j = \{\xi_{(j-1)z+1}, \dots, \xi_{jz}\}$, $j = 1, \dots, 4$, and

$$\Xi_1 \subset \mathcal{B}_\epsilon(+\mathbf{e}_1), \quad \Xi_2 \subset \mathcal{B}_\epsilon(-\mathbf{e}_1), \quad \Xi_3 \subset \mathcal{B}_\epsilon(+\mathbf{e}_2), \quad \Xi_4 \subset \mathcal{B}_\epsilon(-\mathbf{e}_2)$$

and $\xi_{4z+1} = \mathbf{0}$, where $\mathcal{B}_\epsilon(\mathbf{x}) = \{\xi \in \mathbb{R}^d : \|\xi - \mathbf{x}\|_p \leq \epsilon\}$ denotes the ϵ -ball around \mathbf{x} . Here, $\epsilon > 0$ is small enough so that each atom in Ξ_i is closer to $\mathbf{0}$ than to any atom in any of the other sets Ξ_j . The triangle inequality then implies that

$$\begin{aligned} \|\xi_i\|_p &\in [1 - \epsilon, 1 + \epsilon] & \forall \xi_i \in \Xi_1, \\ \|\xi_i - \xi_1\|_p &\geq 2 - 2\epsilon & \forall \xi_i \in \Xi_2, \\ \|\xi_i - \xi_1\|_p &\geq 1 - \epsilon & \forall \xi_i \in \Xi_3 \cup \Xi_4. \end{aligned} \quad (19)$$

Step 2: By construction, we have that

$$\begin{aligned} D_l(\hat{\mathbb{P}}_n, 4) &\leq d_l\left(\hat{\mathbb{P}}_n, \frac{z+1}{4z+1}\delta_{\xi_z} + \frac{z}{4z+1}\delta_{\xi_{2z}} + \frac{z}{4z+1}\delta_{\xi_{3z}} + \frac{z}{4z+1}\delta_{\xi_{4z}}\right) \\ &\leq \left[\frac{1}{4z+1} \left(\sum_{j=1}^4 \sum_{\xi_i \in \Xi_j} \|\xi_i - \xi_{jz}\|_p^l + \|\xi_{4z+1} - \xi_z\|_p^l\right)\right]^{1/l} \\ &\leq \left[\frac{1}{4z+1} \left(\sum_{j=1}^4 \sum_{\xi_i \in \Xi_j} (2\epsilon)^l + (1+\epsilon)^l\right)\right]^{1/l} = \left[\frac{4z2^l\epsilon^l + (1+\epsilon)^l}{4z+1}\right]^{1/l}, \end{aligned}$$

where the first inequality holds because $\xi_z, \xi_{2z}, \xi_{3z}, \xi_{4z} \in \text{supp}(\hat{\mathbb{P}}_n)$, the second inequality holds since moving the atoms in Ξ_j to ξ_{jz} , $j = 1, \dots, 4$, and ξ_{4z+1} to ξ_z represents a feasible transportation plan, and the third inequality is due to (19) and the triangle inequality.

Step 3: We first show that for a sufficiently small $\epsilon > 0$, Dupačová et al.'s algorithm adds $\xi_{4z+1} = \mathbf{0}$ to the reduced set R in the first iteration. We then show that under this selection, the output \mathbb{Q} of Dupačová et al.'s algorithm can be arbitrarily worse than the bound on $D_l(\hat{\mathbb{P}}_n, 4)$ determined in the previous step.

To show the first point, the symmetry inherent in $\text{supp}(\hat{\mathbb{P}}_n)$ implies that it suffices to show that $d_l(\hat{\mathbb{P}}_n, \delta_{\mathbf{0}}) < d_l(\hat{\mathbb{P}}_n, \delta_{\xi_1})$. To this end, we note that

$$d_l^l(\hat{\mathbb{P}}_n, \delta_{\mathbf{0}}) = \frac{1}{4z+1} \sum_{j=1}^4 \sum_{\xi_i \in \Xi_j} \|\xi_i\|_p^l \leq \frac{1}{4z+1} \sum_{j=1}^4 \sum_{\xi_i \in \Xi_j} (1+\epsilon)^l = \frac{4z}{4z+1} (1+\epsilon)^l$$

due to equation (19), while at the same time

$$\begin{aligned} d_l^l(\hat{\mathbb{P}}_n, \delta_{\xi_1}) &= \frac{1}{4z+1} \sum_{i=2}^{4z+1} \|\xi_i - \xi_1\|_p^l \\ &\geq \frac{1}{4z+1} \sum_{i=z+1}^{4z+1} \|\xi_i - \xi_1\|_p^l \geq \frac{z(2-2\epsilon)^l + (2z+1)(1-\epsilon)^l}{4z+1}. \end{aligned}$$

As ϵ tends to 0, we have that

$$\lim_{\epsilon \rightarrow 0} d_l(\hat{\mathbb{P}}_n, \delta_{\mathbf{0}}) \leq \left[\frac{4z}{4z+1}\right]^{1/l} < \left[\frac{z2^l + 2z + 1}{4z+1}\right]^{1/l} \leq \lim_{\epsilon \rightarrow 0} d_l(\hat{\mathbb{P}}_n, \delta_{\xi_1}),$$

where the strict inequality is due to $l \geq 1$. As a consequence, we may conclude that there indeed exists an $\epsilon > 0$ such that Dupačová et al.'s algorithm adds $\xi_{4z+1} = \mathbf{0}$ to the reduced set R in the first iteration.

As for the second point, we note that after adding $\xi_{4z+1} = \mathbf{0}$ to the reduced set R , there must be at least one subset Ξ_j , $j \in \{1, \dots, 4\}$, such that no $\xi_i \in \Xi_j$ is contained in the final reduced set R . Assume w.l.o.g. that this is the case for $j = 1$. We then have

$$d_l(\hat{\mathbb{P}}_n, \mathbb{Q}) \geq \left[\frac{1}{4z+1} \sum_{\xi_i \in \Xi_1} \|\xi_i - \mathbf{0}\|_p \right]^{1/l} \geq \left[\frac{z(1-\epsilon)}{4z+1} \right]^{1/l},$$

and combining this with the result of Step 2, we can conclude that the approximation ratio $d_l(\hat{\mathbb{P}}_n, \mathbb{Q})/D_l(\hat{\mathbb{P}}_n, 4)$ approaches ∞ as $z \rightarrow \infty$ and $z\epsilon^l \rightarrow 0$. \square

We remark that the algorithm of Dupačová et al. can be improved by adding multiple atoms to the reduced set R in Step 2. Nevertheless, a similar argument as in the proof of Theorem 7 shows that the resulting improved algorithm does not allow for a finite upper bound on the approximation ratio either.

4.1.2 k -Means Clustering Algorithm

The k -means clustering algorithm has first been proposed in 1957 for a pulse-code modulation problem (Lloyd 1982), and it has since then become a widely used heuristic for various classes of clustering problems. It aims to partition a set of observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ into m clusters S_1, \dots, S_m such that the intra-cluster sums of squared distances are minimized. By generalizing the algorithm to arbitrary powers and norms, we can adapt the algorithm to our continuous scenario reduction problem as follows.

k -MEANS CLUSTERING ALGORITHM FOR $C_l(\hat{\mathbb{P}}_n, m)$:

1. Initialize the reduced set $R = \{\zeta_1, \dots, \zeta_m\} \subseteq \text{supp}(\hat{\mathbb{P}}_n)$ arbitrarily.
2. Let $\{I_j\} \in \mathfrak{P}(I, m)$ be any partition whose sets I_j , $j \in J$, contain all atoms of $\text{supp}(\hat{\mathbb{P}}_n)$ that are closest to ζ_j (ties may be broken arbitrarily).
3. For each $j \in J$, update $\zeta_j \leftarrow \arg \min \{\sum_{i \in I_j} \|\xi_i - \zeta\|^l : \zeta \in \mathbb{R}^d\}$.
4. Repeat Steps 2 and 3 until the reduced set R no longer changes.

For the empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_{\mathbb{E}}(\mathbb{R}^d, n)$, the algorithm iteratively updates the reduced set R containing the atoms of the reduced distribution \mathbb{Q} through a sequence of assignment (Step 2) and update (Step 3) steps. Step 2 assigns each atom $\xi_i \in \text{supp}(\hat{\mathbb{P}}_n)$ of the empirical distribution to the closest atom in the reduced set, and Step 3 updates each atom in the reduced set so as to minimize the sum of l -th powers of the distances to its assigned atoms from $\text{supp}(\hat{\mathbb{P}}_n)$. After termination, the continuously reduced distribution \mathbb{Q} can be recovered from the reduced set R in the same way as in the previous subsection.

Remark 1 implies that for $l = 2$ and $\|\cdot\| = \|\cdot\|_2$, Step 3 reduces to $\zeta_j \leftarrow \frac{1}{|I_j|} \sum_{i \in I_j} \xi_i$, in which case we recover the classical k -means clustering algorithm. Although the algorithm terminates at a local minimum, Dasgupta (2008) has

shown that for $l = 2$ and $\|\cdot\| = \|\cdot\|_2$, the solution determined by the algorithm can be arbitrarily suboptimal. We now generalize this finding to generic type- l Wasserstein distances induced by arbitrary p -norms.

Theorem 8 *If initialized randomly in Step 1, the approximation ratio of the k -means clustering algorithm is unbounded for every $d, l, p \geq 1$ with significant probability.*

Proof In analogy to the proof of Theorem 7, we construct a specific distribution $\hat{\mathbb{P}}_n$ (Step 1), bound $D_l(\hat{\mathbb{P}}_n, m)$ from above (Step 2) and bound the Wasserstein distance between $\hat{\mathbb{P}}_n$ and the output \mathbb{Q} of the k -means algorithm from below (Step 3).

Step 1: Fix $d, l, p \geq 1$ and $m = 3$, and consider the empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, m)$ with $n = 3z + 1$ for some positive integer z as well as $\text{supp}(\hat{\mathbb{P}}_n) = \Xi_1 \cup \Xi_2 \cup \{\xi_{3z+1}\}$ with $\Xi_1 = \{\xi_1, \dots, \xi_{2z}\}$, $\Xi_2 = \{\xi_{2z+1}, \dots, \xi_{3z}\}$ and

$$\Xi_1 \subset \mathcal{B}_\epsilon(-\mathbf{e}_1), \quad \Xi_2 \subset \mathcal{B}_\epsilon(\mathbf{0}) \quad \text{for } \epsilon \in (0, 1/4), \quad (20)$$

as well as $\xi_{3z+1} = \mathbf{e}_1$, where again $\mathcal{B}_\epsilon(\mathbf{x}) = \{\xi \in \mathbb{R}^d : \|\xi - \mathbf{x}\|_p \leq \epsilon\}$. By construction, the distance between any pair of atoms in Ξ_j is bounded above by $2\epsilon < \frac{1}{2}$, $j = 1, 2$, whereas the distance between two atoms from Ξ_1 and Ξ_2 is bounded from below by $1 - 2\epsilon > \frac{1}{2}$.

Step 2: As similar argument as in the proof of Theorem 7 shows that

$$\begin{aligned} C_l(\hat{\mathbb{P}}_n, 3) &\leq d_l\left(\hat{\mathbb{P}}_n, \frac{2k}{3k+1}\delta_{-\mathbf{e}_1} + \frac{k}{3k+1}\delta_{\mathbf{0}} + \frac{1}{3k+1}\delta_{\mathbf{e}_1}\right) \\ &\leq \left[\frac{1}{3z+1} \left(\sum_{\xi_i \in \Xi_1} \|\xi_i + \mathbf{e}_1\|_p^l + \sum_{\xi_i \in \Xi_2} \|\xi_i\|_p^l \right) \right]^{1/l} \leq \left[\frac{3z\epsilon^l}{3z+1} \right]^{1/l}, \end{aligned}$$

where the last inequality follows from (20).

Step 3: We first show that with significant probability, the algorithm chooses a reduced set R containing two atoms from Ξ_1 and one atom from Ξ_2 in the first step. We then show that under this initialization, the output \mathbb{Q} of the algorithm can be arbitrarily worse than the bound on $C_l(\hat{\mathbb{P}}_n, 3)$ determined above.

In view of the first point, we note that the probability of the reduced set R containing two atoms from Ξ_1 and one atom from Ξ_2 after the first step is $\binom{2z}{2}\binom{z}{1}/\binom{3z+1}{3}$ and approaches 44.44% as $z \rightarrow \infty$. In the following, we thus assume w.l.o.g. that $R = \{\zeta_1, \zeta_2, \zeta_3\}$ with $\zeta_1, \zeta_2 \in \Xi_1$ and $\zeta_3 \in \Xi_2$ after the first step.

As for the second point, we note that Step 2 of the algorithm assigns the atoms $\xi_i \in \Xi_1$ to either ζ_1 or ζ_2 , whereas the atoms $\xi_i \in \Xi_2 \cup \{\xi_{3z+1}\}$ are assigned to ζ_3 . Hence, the update of the reduced set R in the next iteration satisfies $\zeta_1, \zeta_2 \in \mathcal{B}_\epsilon(-\mathbf{e}_1)$, whereas ζ_3 is chosen with respect to the set $\Xi_2 \cup \{\xi_{3z+1}\}$. The algorithm then terminates in the third iteration as the reduced set R no longer changes. We thus find that

$$\begin{aligned} d_l(\hat{\mathbb{P}}_n, \mathbb{Q}) &\geq \left[\frac{1}{3z+1} \sum_{\xi_i \notin \Xi_1} \|\xi_i - \zeta_3\|_p^l \right]^{1/l} \\ &\geq \left[\frac{1}{3z+1} \left(\|\xi_{3z} - \zeta_3\|_p^l + \|\xi_{3z+1} - \zeta_3\|_p^l \right) \right]^{1/l}. \end{aligned}$$

Recall that $\xi_{3z} \in \mathcal{B}_\epsilon(\mathbf{0})$ and $\xi_{3z+1} = \mathbf{e}_1$, which implies that

$$\|\xi_{3z} - \zeta_3\|_p + \|\xi_{3z+1} - \zeta_3\|_p \geq \|\xi_{3z+1} - \xi_{3z}\|_p \geq 1 - \epsilon,$$

and that at least one of the two terms $\|\xi_{3z} - \zeta_3\|_p$ or $\|\xi_{3z+1} - \zeta_3\|_p$ is greater than $\frac{1-\epsilon}{2}$. We thus conclude that

$$d_l(\hat{\mathbb{P}}_n, \mathbb{Q}) \geq \left[\frac{(1-\epsilon)^l}{2^l(3z+1)} \right]^{1/l},$$

which by virtue Step 2 implies that $d_l(\hat{\mathbb{P}}_n, \mathbb{Q})/C_l(\hat{\mathbb{P}}_n, 3) \rightarrow \infty$ as $\epsilon \rightarrow 0$. \square

4.2 Constant-Factor Approximation for the Scenario Reduction Problem

We now propose a simple approximation scheme for the discrete scenario reduction problem under the type-1 Wasserstein distance whose approximation ratio is bounded from above by 5. We also show that this algorithm gives rise to an approximation scheme for the *continuous* scenario reduction problem with an approximation ratio of 10. To our best knowledge, we describe the first constant-factor approximations for the discrete and continuous scenario reduction problems.

Our algorithm follows from the insight that the discrete scenario reduction problem under the type-1 Wasserstein distance is equivalent to the k -median clustering problem. The k -median clustering problem is a variant of the k -means clustering problem described in Section 4.1.2, where the l -th power of the norm terms is dropped (*i.e.*, $l = 1$). In the following, we adapt a well-known local search algorithm (Arya et al 2004) to our discrete scenario reduction problem:

LOCAL SEARCH ALGORITHM FOR $D_l(\hat{\mathbb{P}}_n, m)$:

1. Initialize the reduced set $R \subseteq \text{supp}(\hat{\mathbb{P}}_n)$, $|R| = m$, arbitrarily.
2. Select the next exchange to be applied to the reduced set as

$$(\zeta, \zeta') \in \arg \min \left\{ D_l(\hat{\mathbb{P}}_n, R \cup \{\zeta\} \setminus \{\zeta'\}) : (\zeta, \zeta') \in (\text{supp}(\hat{\mathbb{P}}_n) \setminus R) \times R \right\},$$

and update $R \leftarrow R \cup \{\zeta\} \setminus \{\zeta'\}$ if $D_l(\hat{\mathbb{P}}_n, R \cup \{\zeta\} \setminus \{\zeta'\}) < D_l(\hat{\mathbb{P}}_n, R)$.

3. Repeat Step 2 until no further improvement is possible.

For an empirical distribution $\hat{\mathbb{P}}_n \in \mathcal{P}_E(\mathbb{R}^d, n)$, the algorithm constructs a sequence of reduced sets R containing the atoms of the reduced distribution \mathbb{Q} . In each iteration, Step 2 selects the exchange $R \cup \{\zeta\} \setminus \{\zeta'\}$, $\zeta \in \text{supp}(\hat{\mathbb{P}}_n)$ and $\zeta' \in R$, that maximally reduces the Wasserstein distance $D_l(\hat{\mathbb{P}}_n, R)$. For performance reasons, this ‘best fit’ strategy can also be replaced with a ‘first fit’ strategy which conducts the first exchange $R \cup \{\zeta\} \setminus \{\zeta'\}$ found that leads to a reduction of $D_l(\hat{\mathbb{P}}_n, R)$. After termination, the reduced distribution \mathbb{Q} can be recovered from the reduced set R in the same way as in Section 4.1.1.

It follows from Arya et al (2004) that the above algorithm (with either ‘best fit’ or ‘first fit’) has an approximation ratio of 5 for the discrete scenario reduction problem for all d . We now show that the algorithm also provides solutions to the *continuous* scenario reduction problem with an approximation ratio of at most 10.

Corollary 1 *The problems $D_l(\hat{\mathbb{P}}_n, m)$ and $C_l(\hat{\mathbb{P}}_n, m)$ are related as follows.*

1. *Any approximation algorithm for $D_2(\hat{\mathbb{P}}_n, m)$ under the 2-norm with approximation ratio α gives rise to an approximation algorithm for $C_2(\hat{\mathbb{P}}_n, m)$ under the 2-norm with approximation ratio $\sqrt{2}\alpha$.*
2. *Any approximation algorithm for $D_1(\hat{\mathbb{P}}_n, m)$ under any norm with approximation ratio α gives rise to an approximation algorithm for $C_1(\hat{\mathbb{P}}_n, m)$ under the same norm with approximation ratio 2α .*

Proof The two statements follow directly from Theorems 5 and 6, respectively. \square

As presented, the local search algorithm is not guaranteed to terminate in polynomial time. This can be remedied by a variant of the algorithm that only accepts exchanges $R \cup \{\zeta\} \setminus \{\zeta'\}$ that reduce the Wasserstein distance $D_l(\hat{\mathbb{P}}_n, R)$ by at least $\epsilon/((n-m)m)$ for some constant $\epsilon > 0$. It follows from Arya et al (2004) that for any ϵ , this variant terminates in polynomial time and provides a $(5 + \epsilon)$ -approximation for the discrete scenario reduction problem. The algorithm can also be extended to accommodate multiple swaps in every iteration, which lowers the approximation ratio to $3 + \epsilon$ at the expense of additional computations.

We remark that there is a wealth of algorithms for the k -median problem that can be adapted to the discrete scenario reduction problem. For example, Charikar and Li (2012) present a rounding scheme for the k -median problem that gives rise to a polynomial-time algorithm for $D_1(\hat{\mathbb{P}}_n, R)$ with an approximation ratio of 3.25. Likewise, Kanungo et al (2004) propose a local search algorithm for the k -median problem that gives rise to a polynomial-time algorithm for $D_2(\hat{\mathbb{P}}_n, R)$ under the 2-norm with an approximation ratio of $9 + \epsilon$. In both cases, Corollary 1 allows us to extend these guarantees to the corresponding versions of the continuous scenario reduction problem.

4.3 Mixed-Integer Reformulations of the Discrete and Continuous Scenario Reduction Problems

We first review a well-known mixed-integer linear programming (MILP) reformulation of the discrete scenario reduction problem $D_l(\hat{\mathbb{P}}_n, m)$:

Theorem 9 *The discrete scenario reduction problem can be formulated as the MILP*

$$\begin{aligned} D_l^l(\hat{\mathbb{P}}_n, m) &= \min_{\mathbf{\Pi}, \boldsymbol{\lambda}} \quad \frac{1}{n} \langle \mathbf{\Pi}, \mathbf{D} \rangle \\ \text{s.t.} \quad &\mathbf{\Pi} \mathbf{e} = \mathbf{e}, \quad \mathbf{\Pi} \leq \mathbf{e} \boldsymbol{\lambda}^\top, \quad \boldsymbol{\lambda}^\top \mathbf{e} = m \\ &\mathbf{\Pi} \in \mathbb{R}_+^{n \times n}, \quad \boldsymbol{\lambda} \in \{0, 1\}^n, \end{aligned} \tag{21}$$

with $\mathbf{D} \in \mathbb{S}^n$ and $d_{ij} = \|\boldsymbol{\xi}_i - \boldsymbol{\xi}_j\|^l$ encoding the distances among the atoms in $\text{supp}(\hat{\mathbb{P}}_n)$.

Proof See e.g. Heitsch and Römisch (2003). \square

In problem (21), the decision variable π_{ij} determines how much of the probability mass of atom $\boldsymbol{\xi}_i$ in $\hat{\mathbb{P}}_n$ is shifted to the atom $\boldsymbol{\xi}_j$ in the reduced distribution \mathbb{Q} , whereas the decision variable λ_j determines whether the atom $\boldsymbol{\xi}_j \in \text{supp}(\hat{\mathbb{P}}_n)$ is contained in the support of \mathbb{Q} . A solution $(\mathbf{\Pi}^*, \boldsymbol{\lambda}^*)$ to problem (21) allows us to

recover the reduced distribution via $\mathbb{Q} = \frac{1}{n} \sum_{j=1}^n \mathbf{e}^\top \mathbf{\Pi}^* \mathbf{e}_j \cdot \delta_{\boldsymbol{\xi}_j}$. Problem (21) has n binary and n^2 continuous variables as well as $n^2 + n + 1$ constraints.

We now consider the *continuous* scenario reduction problem. Due to its bilinear objective function, which involves products of transportation weights π_{ij} and the distances $\|\boldsymbol{\xi}_i - \boldsymbol{\zeta}_j\|^l$ containing the continuous decision variables $\boldsymbol{\zeta}_j$, this problem may not appear to be amenable to a reformulation as a mixed-integer convex optimization problem. We now show that such a reformulation indeed exists.

Theorem 10 *The continuous scenario reduction problem can be formulated as the mixed-integer convex optimization problem*

$$\begin{aligned} C_l^l(\hat{\mathbb{P}}_n, m) = \min_{\mathbf{\Pi}, \mathbf{c}, \{\boldsymbol{\zeta}_j\}} \quad & \frac{1}{n} \mathbf{e}^\top \mathbf{c} \\ \text{s.t.} \quad & \mathbf{\Pi} \mathbf{e} = \mathbf{e} \\ & \|\boldsymbol{\xi}_i - \boldsymbol{\zeta}_j\|^l \leq c_i + M(1 - \pi_{ij}) \quad \forall i \in I, \forall j \in J \\ & \mathbf{\Pi} \in \{0, 1\}_+^{n \times m}, \quad \mathbf{c} \in \mathbb{R}_+^n, \quad \boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_m \in \mathbb{R}^d, \end{aligned} \quad (22)$$

where $M = \max_{i,j \in I} \|\boldsymbol{\xi}_i - \boldsymbol{\xi}_j\|^l$ denotes the diameter of the support of $\hat{\mathbb{P}}_n$.

In problem (22), the decision variable π_{ij} determines whether or not the probability mass of atom $\boldsymbol{\xi}_i$ in the empirical distribution $\hat{\mathbb{P}}_n$ is shifted to the atom $\boldsymbol{\zeta}_j$ in the reduced distribution \mathbb{Q} , whereas the decision variable c_i records the cost of moving the atom $\boldsymbol{\xi}_i$ under the transportation plan $\mathbf{\Pi}$.

Proof of Theorem 10 We prove the statement by showing that optimal solutions to problem (22) correspond to feasible solutions in problem (2) with the same objective function value in their respective problems and vice versa.

Fix a minimizer $(\mathbf{\Pi}^*, \mathbf{c}^*, \boldsymbol{\zeta}_1^*, \dots, \boldsymbol{\zeta}_m^*)$ to problem (22), which corresponds to a feasible solution $(\{I_j\}, \boldsymbol{\zeta}_1^*, \dots, \boldsymbol{\zeta}_m^*)$ in problem (2) if we set $I_j = \{i \in I : \pi_{ij}^* = 1\}$ for all $j \in J$. Note that $c_i^* = \|\boldsymbol{\xi}_i - \boldsymbol{\zeta}_j^*\|^l$ for $j \in J$ and $i \in I_j$. Thus, both solutions adopt the same objective value.

Conversely, fix a minimizer $(\{I_j^*\}, \boldsymbol{\zeta}_1^*, \dots, \boldsymbol{\zeta}_m^*)$ to problem (2). This solution corresponds to a feasible solution $(\mathbf{\Pi}, \mathbf{c}, \boldsymbol{\zeta}_1^*, \dots, \boldsymbol{\zeta}_m^*)$ to problem (22) if we set $\pi_{ij} = 1$ if $i \in I_j^*$ and $\pi_{ij} = 0$ otherwise for all $j \in J$, as well as $c_i = \|\boldsymbol{\xi}_i - \boldsymbol{\zeta}_j^*\|^l$ for all $i \in I_j$ and $j \in J$. By construction, both solutions adopt the same objective value. \square

A solution $(\mathbf{\Pi}^*, \mathbf{c}^*, \boldsymbol{\zeta}_1^*, \dots, \boldsymbol{\zeta}_m^*)$ to problem (22) allows us to recover the reduced distribution via $\mathbb{Q} = \frac{1}{n} \sum_{j=1}^m \mathbf{e}^\top \mathbf{\Pi}^* \mathbf{e}_j \cdot \delta_{\boldsymbol{\zeta}_j^*}$. Problem (22) has nm binary and $n + md$ continuous variables as well as $nm + n$ constraints. We now show that (22) typically reduces to an MILP or a mixed-integer second-order cone program (MISOCP).

Proposition 8 *For the type-1 Wasserstein distance induced by $\|\cdot\|_1$ or $\|\cdot\|_\infty$, problem (22) reduces to an MILP. For any type- l Wasserstein distance induced by $\|\cdot\|_p$, where $l \geq 1$ and $p \geq 1$ are rational numbers, problem (22) reduces to an MISOCP.*

Proof In view of the first statement, we note that $\|\boldsymbol{\xi}_i - \boldsymbol{\zeta}_j\|_1 \leq c_i + M(1 - \pi_{ij})$ is satisfied if and only if there is $\boldsymbol{\phi}_{ij} \in \mathbb{R}^d$ such that

$$\boldsymbol{\phi}_{ij} \geq \boldsymbol{\xi}_i - \boldsymbol{\zeta}_j, \quad \boldsymbol{\phi}_{ij} \geq \boldsymbol{\zeta}_j - \boldsymbol{\xi}_i \quad \text{and} \quad \mathbf{e}^\top \boldsymbol{\phi}_{ij} \leq c_i + M(1 - \pi_{ij}).$$

Likewise, $\|\xi_i - \zeta_j\|_\infty \leq c_i + M(1 - \pi_{ij})$ holds if and only if there is $\phi_{ij} \in \mathbb{R}$ with

$$\phi_{ij}\mathbf{e} \geq \xi_i - \zeta_j, \quad \phi_{ij}\mathbf{e} \geq \zeta_j - \xi_i \quad \text{and} \quad \phi_{ij} \leq c_i + M(1 - \pi_{ij}).$$

As for the second statement, we note that $\|\xi_i - \zeta_j\|_p^l \leq c_i + M(1 - \pi_{ij})$ is satisfied if and only if there is $\phi_{ij} \in \mathbb{R}$ such that

$$\phi_{ij} \geq \|\xi_i - \zeta_j\|_p \quad \text{and} \quad \phi_{ij}^l \leq c_i + M(1 - \pi_{ij}).$$

For rational $l, p \geq 1$, both inequalities can be expressed through finitely many second-order cone constraints (Alizadeh and Goldfarb 2003, Section 2.3). \square

5 Numerical Experiment: Color Quantization

Color quantization aims to reduce the color palette of a digital image without compromising its visual appearance. In the standard RGB24 model colors are encoded by vectors of the form $(r, g, b) \in \{0, 1, \dots, 255\}^3$. This means that the RGB24 model can represent a vast number of 16,777,216 distinct colors. Consequently, color quantization serves primarily as a lossy image compression method.

In the following we interpret the color quantization problem as a discrete scenario reduction problem using the type-1 Wasserstein distance induced by the 1-norm on \mathbb{R}^3 . Thus, we can solve color quantization problems via Dupačová’s greedy heuristic, the local search algorithm or the exact MILP reformulation (21). In our experiment we aim to compress all 24 pictures from the Kodak Lossless True Color Image Suite (<http://r0k.us/graphics/kodak/>) to $m = 2^1, \dots, 2^9$ colors. As the MILP reformulation scales poorly with n , we first reduce each image to $n \lesssim 1,024$ colors using the Linux command “convert -colors”, which is distributed through ImageMagick (<https://www.imagemagick.org>). We henceforth refer to the resulting 1,024-color images as the originals.

In all experiments we use an efficient variant of Dupačová’s algorithm due to Heitsch and Römischi (2003) (DPCV), and we initialize the local search algorithm either with the color palette obtained from Dupačová’s algorithm (LOC-1) or naïvely with the m most frequent colors of the original image (LOC-2). The MILP (21) is solved with GUROBI 7.0.1 (MILP). All algorithms are implemented in C++, and all experiments are executed on a 3.40GHz i7 CPU machine with 16GB RAM. We report the average and the worst-case runtimes in Table 1. Note that DPCV, LOC-1 and LOC-2 all terminate in less than 14 seconds across all instances, while MILP requires substantially more time (the maximum runtime was set to ten hours). Moreover, warmstarting the local search algorithm with the color palette obtained from DPCV can significantly reduce the runtimes.

	DPCV	LOC-1	LOC-2	MILP
Average (secs)	2.16	2.60	3.59	1,349.71
Worst-case (secs)	8.21	9.89	13.69	36,120.99

Table 1: Runtimes of different methods for discrete scenario reduction.

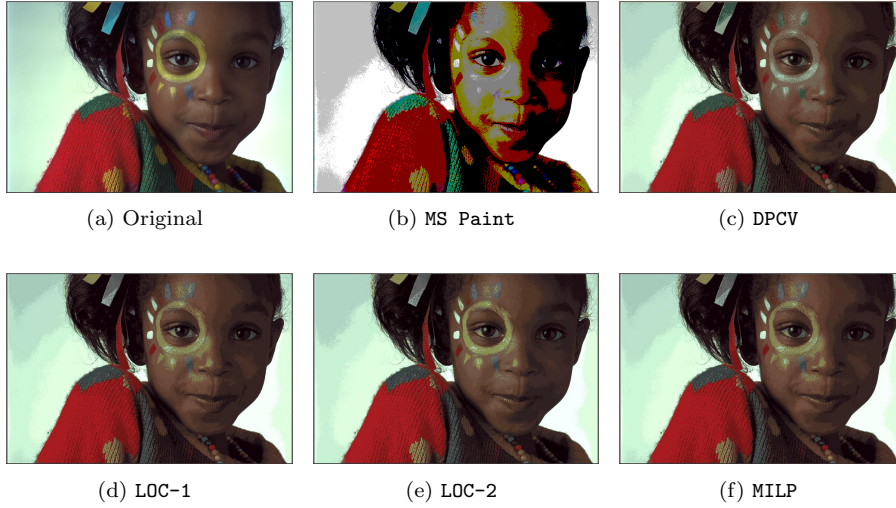


Fig. 3: Outputs of different color quantization algorithms for image “kodim15.png”.

As an example, Figure 3 shows the image “kodim15.png” as well as the results of different color quantization algorithms for $m = 16$. While the outputs of LOC-1, LOC-2 and MILP are almost indistinguishable, the output of DPCV has ostensible deficiencies (*e.g.*, it misrepresents the yellow color around the subject’s eye). For comparison, we also show the output of the color quantization routine in Microsoft Paint (MS Paint). Figure 4 visualizes the optimality gaps of DPCV, LOC-1 and LOC-2 relative to MILP (*i.e.*, their respective approximation ratio $- 1$). Our experiment suggests that the local search algorithm is competitive with MILP in terms of output quality but at significantly reduced runtimes. Moreover, the local search algorithm LOC-1 warmstarted with the color palette obtained from DPCV is guaranteed to outperform DPCV in terms of optimality gaps.

Acknowledgements This research was funded by the SNSF grant BSCGI0_157733 and the EPSRC grants EP/M028240/1 and EP/M027856/1.

References

- Alizadeh F, Goldfarb D (2003) Second-order cone programming. *Mathematical Programming* 95(1):3–51
- Aloise D, Deshpande A, Hansen P, Popat P (2009) NP-hardness of Euclidean sum-of-squares clustering. *Machine Learning* 75(2):245–248
- Arya V, Garg N, Khandekar R, Meyerson A, Munagala K, Pandit V (2004) Local search heuristics for k -median and facility location problems. *SIAM Journal on Computing* 33(3):544–562
- Charikar M, Li S (2012) A dependent LP-rounding approach for the k -median problem. In: *Proceedings of the 39th International Colloquium Conference on Automata, Languages, and Programming*, pp 194–205

- Dasgupta S (2008) CSE 291: Topics in unsupervised learning. URL <http://cseweb.ucsd.edu/~dasgupta/291-unsup/>
- Dupačová J (1990) Stability and sensitivity-analysis for stochastic programming. *Annals of Operations Research* 27(1):115–142
- Dupačová J, Gröwe-Kuska N, Römisch W (2003) Scenario reduction in stochastic programming: an approach using probability metrics. *Mathematical Programming* 95(3):493–511
- Gao R, Kleywegt A (2016) Distributionally robust stochastic optimization with Wasserstein distance. *arXiv* #1604.02199
- Graf S, Luschgy H (2000) *Foundations of Quantization for Probability Distributions*. Springer
- Gray RM (2006) *Toeplitz and Circulant Matrices: A Review*. Now Publishers
- Hanasusanto G, Kuhn D (2016) Conic programming reformulations of two-stage distributionally robust linear programs over Wasserstein balls. *arXiv* #1609.07505
- Heitsch H, Römisch W (2003) Scenario reduction algorithms in stochastic programming. *Computational Optimization and Applications* 24(2):187–206
- Heitsch H, Römisch W (2007) A note on scenario reduction for two-stage stochastic programs. *Operations Research Letters* 35(6):731 – 738
- Hochreiter R, Pflug G (2007) Financial scenario generation for stochastic multi-stage decision processes as facility location problems. *Annals of Operations Research* 152(1):257–272
- Hopcroft J, Kannan R (2012) *Computer science theory for the information age*. URL <https://www.cs.cmu.edu/~venkatg/teaching/CStheory-infoage/>
- Kanungo T, Mount D, Netanyahu N, Piatko C, Silverman R, Wu A (2004) A local search approximation algorithm for k -means clustering. *Computational Geometry* 28(2):89–112
- Kariv O, Hakimi SL (1979) An algorithmic approach to network location problems. ii: The p -medians. *SIAM Journal on Applied Mathematics* 37(3):539–560
- Lloyd S (1982) Least squares quantization in PCM. *IEEE Transactions on Information Theory* 28(2):129–137
- Mahajan M, Nimbhorkar P, Varadarajan K (2009) The planar k -means problem is NP-hard. In: *Proceedings of the 3rd International Workshop on Algorithms and Computation*, pp 274–285
- Mohajerin Esfahani P, Kuhn D (2015) Data-driven distributionally robust optimization using the Wasserstein metric: Performance guarantees and tractable reformulations. *arXiv* #1505.05116
- Morales JM, Pineda S, Conejo AJ, Carrion M (2009) Scenario reduction for futures market trading in electricity markets. *IEEE Transactions on Power Systems* 24(2):878–888
- Pflug G (2001) Scenario tree generation for multiperiod financial optimization by optimal discretization. *Mathematical Programming* 89(2):251–271
- Pflug G, Pichler A (2011) Approximations for probability distributions and stochastic optimization problems. In: Bertocchi M, Consigli G, Dempster MAH (eds) *Stochastic Optimization Methods in Finance and Energy: New Financial Products and Energy Market Strategies*, Springer, pp 343–387
- Römisch W (2003) Stability of stochastic programming problems. In: Ruszczyński A, Shapiro A (eds) *Stochastic Programming*, Elsevier, pp 483–554
- Römisch W, Vigerske S (2010) Recent progress in two-stage mixed-integer stochastic programming with applications to power production planning. In: Pardalos P, Rebennack S, Pereira M, Iliadis N (eds) *Handbook of Power Systems I*, Springer, pp 177–208
- Steele M (2004) *The Cauchy-Schwarz Master Class: An Introduction to the Art of Mathematical Inequalities*. Cambridge University Press
- Zhao C, Guan Y (2015) Data-driven risk-averse stochastic optimization with Wasserstein metric. Available on Optimization Online

Appendix: Auxiliary Results

The proofs of Theorem 2 relies on the following two lemmas.

Lemma 1 *The semidefinite program (6) admits an optimal solution (τ, \mathbf{S}) with $\mathbf{S} = \alpha \mathbb{I} + \beta \mathbf{1}\mathbf{1}^\top$ for some $\alpha, \beta \in \mathbb{R}$.*

Proof Let (τ, \mathbf{S}^*) be any optimal solution to (6), which exists because (6) has a continuous objective function and a compact feasible set, and denote by \mathfrak{S} the set of all permutations of I . For any $\sigma \in \mathfrak{S}$, the permuted solution $(\tau, \mathbf{S}^\sigma)$, with $s_{ij}^\sigma = s_{\sigma(i)\sigma(j)}^*$ is also optimal in (6). Note first that $(\tau, \mathbf{S}^\sigma)$ is feasible in (6) because

$$\begin{aligned} \tau &\leq \sum_{j \in J} \frac{1}{|I_j|^2} \sum_{i \in I_j} \left(|I_j|^2 s_{ii}^\sigma - 2|I_j| \sum_{k \in I_j} s_{ik}^\sigma + \sum_{k \in I_j} s_{kk}^\sigma + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} s_{kk'}^\sigma \right) \\ \iff \tau &\leq \sum_{j \in J} \frac{1}{|I_j^\sigma|^2} \sum_{i \in I_j^\sigma} \left(|I_j^\sigma|^2 s_{ii}^* - 2|I_j^\sigma| \sum_{k \in I_j} s_{ik}^* + \sum_{k \in I_j} s_{kk}^* + \sum_{\substack{k, k' \in I_j \\ k \neq k'}} s_{kk'}^* \right), \end{aligned}$$

where the index sets $I_j^\sigma = \{\sigma(i) : i \in I_j\}$ for $j \in J$ form an m -set partition from within $\mathfrak{P}(I, m)$, and because $\mathbf{S}^\sigma \succeq \mathbf{0}$ and $s_{ii}^\sigma = s_{\sigma(i)\sigma(i)}^* \leq 1$ for all $i \in I$ by construction. Moreover, it is clear that $(\tau, \mathbf{S}^\sigma)$ and (τ, \mathbf{S}^*) share the same objective value in (6). Thus, $(\tau, \mathbf{S}^\sigma)$ is optimal in (6) for every $\sigma \in \mathfrak{S}$.

The convexity of problem (6) implies that (τ, \mathbf{S}) with $\mathbf{S} = \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}} \mathbf{S}^\sigma$ is also optimal in (6). The claim follows by noting that \mathbf{S} is invariant under permutations of the coordinates and thus representable as $\alpha \mathbb{I} + \beta \mathbf{1}\mathbf{1}^\top$ for some $\alpha, \beta \in \mathbb{R}$. \square

Lemma 2 *For $\alpha, \beta \in \mathbb{R}$ the eigenvalues of $\mathbf{S} = \alpha \mathbb{I} + \beta \mathbf{1}\mathbf{1}^\top \in \mathbb{S}^n$ are given by $\alpha + n\beta$ (with multiplicity 1) and α (with multiplicity $n - 1$).*

Proof Note that \mathbf{S} is a circulant matrix, meaning that each of its rows coincides with the preceding row rotated by one element to the right. Thus, the eigenvalues of \mathbf{S} are given by $\alpha + \beta(1 + \rho_j^1 + \dots + \rho_j^{n-1})$, $j = 0, \dots, n-1$, where $\rho_j = e^{2\pi i j/n}$ and i denotes the imaginary unit; see *e.g.* Gray (2006). For $j = 0$ we then obtain the eigenvalue $\alpha + n\beta$, and for $j = 1, \dots, n-1$ we obtain the other $n-1$ eigenvalues, all of which equal α because $\sum_{k=0}^{n-1} e^{2\pi i j k/n} = (1 - e^{2\pi i j})/(1 - e^{2\pi i j/n}) = 0$. \square

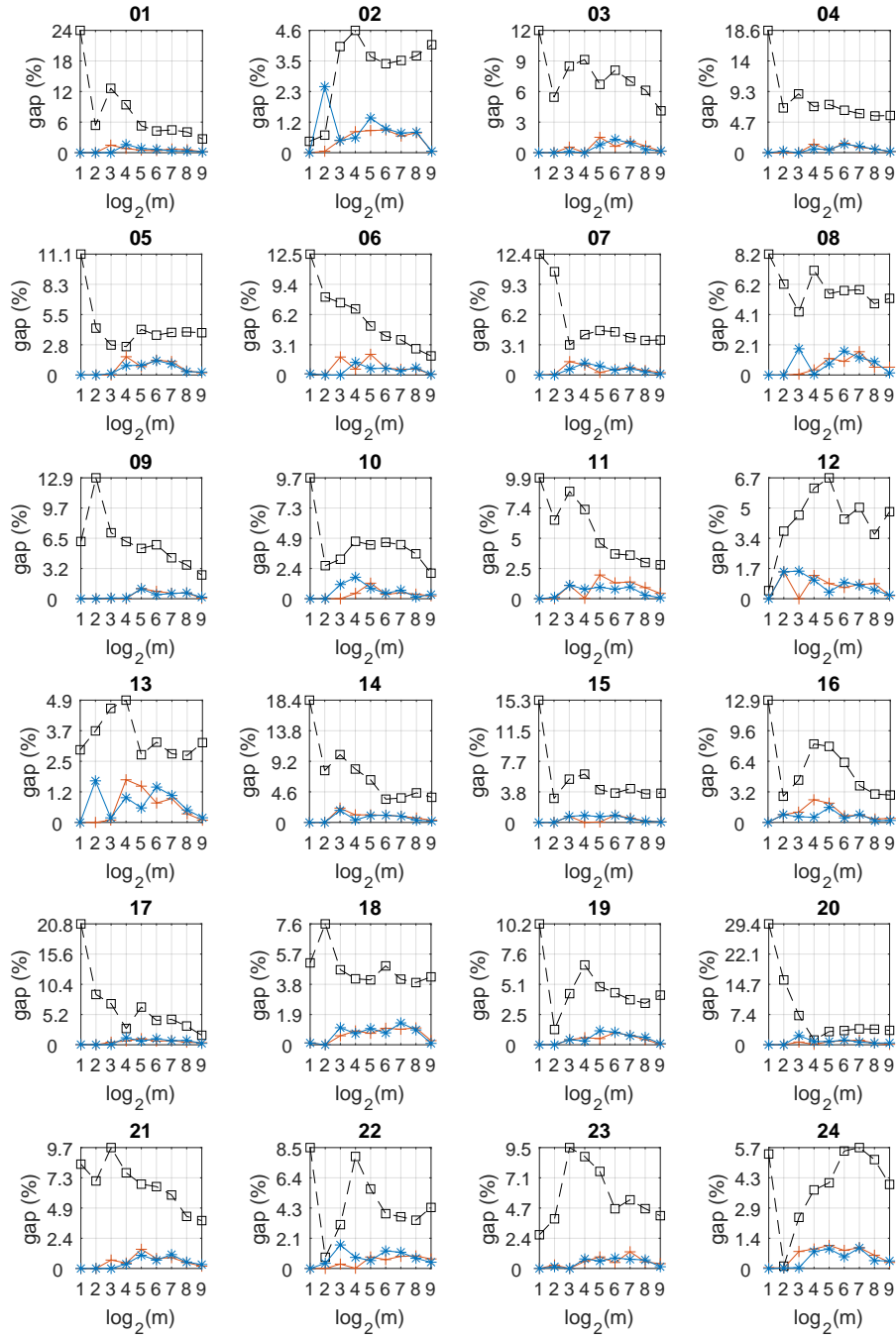


Fig. 4: Optimality gaps of DPCV (dashed line with box), LOC-1 (solid line with plus) and LOC-2 (solid line with star) relative to MILP.